

INVESTIGAÇÃO OPERACIONAL

Junho 1989

Número 1

Volume 9

Publicação Científica da



Associação Portuguesa para o Desenvolvimento
da Investigação Operacional.

INVESTIGAÇÃO OPERACIONAL

Propriedade:

APDIO — Associação Portuguesa para o Desenvolvimento
da Investigação Operacional

ESTATUTO EDITORIAL

«Investigação Operacional», órgão oficial da APDIO cobre uma larga gama de assuntos reflectindo assim a grande diversidade de profissões e interesses dos sócios da Associação, bem como as muitas áreas de aplicação da I. O. O seu objectivo primordial é promover a aplicação do método e técnicas da I. O. aos problemas da Sociedade Portuguesa.

A publicação acolhe contribuições nos campos da metodologia, técnicas, e áreas de aplicação e software de I. O. sendo no entanto dada prioridade a bons casos de estudo de carácter eminentemente prático.

Distribuição gratuita aos sócios da APDIO

INVESTIGAÇÃO OPERACIONAL

volume 9 - nº 1 - Junho 1989

Publicação semestral

Direcção : J.M. Pinto Paixão
(Fac. Ciências - Universidade de Lisboa)
Joaquim J. Júdice
(Fac. Ciências e Tecnologia - Univ. de Coimbra)

Comissão Editorial

Mordecai Avriel	(Israel)	Nelson Maculan	(UFRJ - Brasil)
Paulo Bárcia	(FE - Univ. Nova Lisboa)	A. Simões Monteiro	(NORMA)
João A. Branco	(IST - Univ. Técn. Lisboa)	Mohamed Najim	(ENSIAS - Argélia)
João Clímaco	(INESC/Univ.Coimbra)	J. Manuel Oliveira	(EFACEC)
Josep Casanovas	(UPC - Espanha)	Fernando Pacheco	(Univ. Católica)
J. Dias Coelho	(FE - Univ. Nova Lisboa)	A. Gouvêa Portela	(IST- Univ.Técn. Lisboa)
Nuno Crato	(NORMA - Açores)	M. Baptista Rodrigues	(Partex - CPS)
J.A.Romão Eusébio	(CIMPOR)	A. Guimarães Rodrigues	(Univ. Minho)
A. Sousa Ferraria	(Petrogal)	Bernard Roy	(LAMSADE- França)
D. V. Gokhale	(Estados Unidos)	C. Moreira da Silva	(FE - Univ. Porto)
J. Borges Gouveia	(FE - Univ. Porto)	L. Valadares Tavares	(IST- Univ.Técn. Lisboa)
R. Campos Guimarães	(FE - Univ. Porto)	Isabel H. Themido	(IST- Univ.Técn. Lisboa)
Masao Iri	(TU - Japão)	B. Calafate Vasconcelos	(FE - Univ. Porto)
Joaquim J. Júdice	(FC - Univ. Coimbra)	José M. Viegas	(IST- Univ.Técn. Lisboa)
A. Rinnoy Kan	(EU - Holanda)	Andres Weintraub	(UC - Chile)

A Revista "INVESTIGAÇÃO OPERACIONAL" está registada na Secretaria de Estado da Comunicação Social sob o nº108335.

Esta Revista é distribuída gratuitamente aos sócios da APDIO. As informações sobre inscrições na Associação, assim como a correspondência para a Revista devem ser enviadas para a sede da APDIO - Associação Portuguesa para o Desenvolvimento da Investigação Operacional - CESUR, Instituto Superior Técnico, Av. Rovisco Pais, 1000 Lisboa.

Este Volume foi subsidiado por :
Instituto Nacional de Investigação Científica (INIC)
Junta Nacional de Investigação Científica e Tecnológica (JNICT)
Fundação Calouste Gulbenkian.
Para efeitos de dactilografia e composição, foram utilizados equipamentos gentilmente postos à disposição pelo CERUL (DEIOC-Faculdade de Ciências de Lisboa).

Assinatura : 3000\$00

PROBLEMA BÁSICO DE DISTRIBUIÇÃO E SUAS EXTENSÕES - UMA REVISÃO BIBLIOGRÁFICA -

Maria Teresa Almeida
Instituto Superior de Economia e Gestão (UTL)
Rua Miguel Lupi, 20, 1200 Lisboa

Resumo: O problema básico de distribuição consiste na determinação da configuração óptima de um conjunto de rotas a serem percorridas por uma frota de veículos para servirem um conjunto de clientes cujas localizações e procuras são conhecidas, minimizando os custos de distribuição e respeitando as restrições de capacidade.

Este tipo de problemas encontra-se em numerosas situações de planeamento de transporte de pessoas e bens, bem como em situações de planeamento de determinadas tarefas. Ele pode ser generalizado à existência de numerosas outras restrições e/ou condições de planeamento, dando origem a uma vasta classe de problemas, geralmente designados por problemas de distribuição.

Neste trabalho é feita uma revisão bibliográfica dos principais métodos exactos e aproximados de resolução de problemas de distribuição.

1. Introdução

A designação "Problemas de Distribuição", tradução das designações "Vehicle Routing", "Vehicle Scheduling", "Truck Dispatching" e "Delivery Problem" encontradas na literatura especializada de língua inglesa, engloba uma enorme variedade de problemas. Se em muitos deles se trata efectivamente de distribuir (e/ou recolher) produtos (derivados do petróleo, correspondência postal, produtos alimentares, lixo doméstico e industrial, etc.), noutros trata-se do transporte de pessoas (transportes especiais para deficientes, transportes escolares, etc.) e noutros ainda trata-se de situações em que não há bens ou pessoas a transportar, mas antes, determinadas operações a efectuar (inspecção de redes ou condutas, selecção de produtos num armazém, etc.) as quais podem ser modelizadas como se se tratasse de um problema do tipo dos anteriores. A enorme variedade de situações englobadas conduz, naturalmente, a uma multiplicidade de modelos.

Pode, no entanto, considerar-se que na base de todos eles está o problema, neste trabalho designado por problema básico de distribuição, da determinação do percurso de comprimento mínimo a efectuar por uma frota de veículos estacionada num armazém central para servir um conjunto de clientes. Conhecidas as capacidades dos veículos, as localizações geográficas dos clientes e do armazém central e as procuras, trata-se de determinar as configurações das rotas que correspondem ao menor comprimento total. Para isso é necessário, por um lado, proceder à afectação dos clientes às rotas, respeitando as restrições de capacidade e, por outro, determinar a ordem pela qual deverão ser visitados os clientes incluídos numa mesma rota.

O problema básico de distribuição pode ser considerado uma generalização do problema do caixeiro viajante múltiplo, Bodin et al (1983), obtida pela introdução de restrições na carga total a transportar em cada rota. Este, por sua vez, pode ser visto como uma extensão do clássico problema do caixeiro viajante, Christofides (1979), em que se admitem repetidas passagens pelo armazém central.

Todos os problemas já referidos pertencem à classe dos problemas NP-completos, Lenstra e Rinnooy Kan (1981), não se conhecendo por isso, para nenhum deles, algoritmos exactos de complexidade computacional. Contudo, os problemas de distribuição têm-se mostrado na prática comparativamente ainda mais difíceis. A dimensão dos problemas de distribuição para os quais é possível encontrar solução exacta está muito aquém não só da dimensão da maioria das aplicações reais, como da dimensão já alcançada para o problema do caixeiro viajante, Balas e Christofides (1981), Crowder e Padberg (1980) e Padberg e Rinaldi (1986). Laporte e Norbert (1987) indicam como dimensão máxima já resolvida (excepto para certas condições de tipo muito particular) cerca de 30 clientes. Lucena Filho (1986) na sua tese, resolve problemas com 45 clientes, resultado que é apresentado como correspondendo aproximadamente à duplicação da dimensão máxima resolvida até então.

Na Secção 2 são apresentadas as principais extensões do problema básico de distribuição. Na Secção 3 são apresentadas formalizações alternativas do problema básico de distribuição e métodos exactos de optimização nelas baseados. Na Secção 4 são apresentados métodos heurísticos para a obtenção de soluções aproximadas.

2. Problema Básico de Distribuição e suas Extensões

No problema básico de distribuição, como já foi referido, considera-se um conjunto de n localizações geográficas, $N = \{1, 2, \dots, n\}$, a primeira representando um armazém central onde se encontra estacionada uma frota de veículos e as restantes representando os clientes. Para cada cliente, além da sua localização, é também conhecida a sua procura q_i , $i = 2, \dots, n$, relativa ao bem em consideração. Para cada veículo é conhecida a sua capacidade Q_k , $k = 1, \dots, M$. Trata-se de determinar as configurações das rotas a efectuar pelos veículos, de modo a que cada cliente seja servido por um e um só veículo, minimizando-se o comprimento (ou custo) do percurso total.

Esta situação pode comportar diversas restrições adicionais. De entre as mais frequentes podem salientar-se as seguintes:

(I) Restrições Temporais

Os veículos (ou, mais exactamente, os seus condutores) só podem operar durante intervalos de tempo de duração limitada. Assim, cada rota não pode ter uma duração superior a um determinado valor pré-fixado. Na avaliação da duração da rota, dependendo das situações, podem ser apenas incluídos os tempos de trajecto, ou ser também incluídos os tempos de carregamento do veículo no armazém central e os tempos de descarga e entrega das encomendas nos clientes.

Os clientes só aceitam as encomendas durante certos intervalos de tempo (designados na literatura inglesa por "time-windows"). Neste caso, podem ainda distinguir-se as situações em que se admite que o veículo chegue aos clientes fora dos intervalos temporais admissíveis e fique estacionado à espera que eles se iniciem e as situações em que devido à inexistência de local de estacionamento, o veículo só pode chegar aos clientes durante os intervalos temporais admissíveis para as descargas. No primeiro caso os tempos de espera devem também ser considerados na avaliação da duração temporal da respectiva rota. Solomon e Desrosiers (1987) fazem uma revisão bibliográfica sobre estes problemas.

(II) Restrições de Precedência

Nalguns casos há restrições de precedência entre clientes, não podendo determinados clientes ser visitados sem que outros o tenham já sido. Esta situação é particularmente frequente nos modelos em que há simultaneamente recolha e distribuição, pois em muitos casos, determinadas recolhas têm de preceder algumas distribuições e, por outro lado, a sequência de recolhas e distribuição, tem de ser compatível com as capacidades dos veículos.

(III) Restrições na Frota

O número de veículos que compõem a frota pode ser conhecido a priori tendo nesse caso que impôr-se a restrição adicional de que o número de rotas a gerar não pode ultrapassar o número de veículos disponível.

O problema básico de distribuição pode também ser generalizado, sendo as suas generalizações mais frequentes as seguintes:

(I) Múltiplos Depósitos

Se existir mais do que um armazém para estacionamento dos veículos, mas cada veículo estiver afectado a priori a um dos armazéns (por razões de organização regional ou outras), o problema é decomponível em tantos problemas com um único armazém, quantos os armazéns existentes. Se, no entanto, se admitir que um veículo pode iniciar a sua rota num armazém e terminá-la noutro, além da afectação dos clientes às rotas é ainda necessário proceder à afectação das partidas e chegadas de cada veículo a um armazém, tendo em atenção a capacidade de estacionamento de cada um deles.

(II) Frota Não Homogénea

Nalguns casos a frota dos veículos não é homogénea, não só porque nem todos os veículos têm a mesma capacidade, mas porque nem todos têm as mesmas condições de transporte (compartimento frigorífico, por exemplo, no caso de bens alimentares). Sendo assim é necessário proceder à afectação dos veículos às rotas geradas tendo em consideração as respectivas características.

(III) Múltiplos objectivos

Para além do objectivo já referido da minimização do comprimento (ou custo) total, podem estabelecer-se outros objectivos. Quando o número de veículos na frota não é pré-determinado (porque se admite a compra ou aluguer de novos veículos, se necessário) além da minimização dos custos de transporte, designados geralmente por custos variáveis, tem de incluir-se também na função objectivo uma parcela relativa aos custos de compra ou aluguer dos veículos utilizados, geralmente designados por custos fixos. Podem ainda incluir-se na função objectivo penalizações relativas a clientes não servidos quando com as restrições não é imposto que cada cliente seja visitado por um e um só veículo, ou ponderadores dos clientes que favoreçam os que sejam considerados de maior importância.

Para além destas extensões do problema básico de distribuição, merecem ainda particular destaque, aquelas que se referem a tipos de problemas bem caracterizados como os de transporte escolar, Desrosiers et al (1986), Graham e Nuttle (1986) e Gavish e Graves (1979), e os transportes partilhados de passageiros ("dial-a-ride" na literatura inglesa), Desrosiers, Dumas e Soumis (1984, 1986), Bodin e Sexton (1986), Gavish e Srikanth (1979) e Gavish e Graves (1979). No caso dos transportes escolares são particularmente relevantes as condições referentes a normas de segurança. No caso dos transportes partilhados de passageiros cada cliente corresponde a um par de localizações, com relação de precedência entre si, referentes aos pontos de início e fim da viagem a efectuar.

Por último, refira-se ainda que, por vezes, o planeamento das rotas dos veículos é feito com o objectivo de o repetir periodicamente ao longo de um horizonte temporal, sendo neste caso utilizada a designação de "rotas fixas", Beasley (1984).

3. Algoritmos Exactos

Dada a complexidade computacional do problema básico de distribuição, todos os algoritmos exactos desenvolvidos para a sua resolução, têm de recorrer às técnicas de programação inteira, em particular, à pesquisa em árvore.

Estes algoritmos baseiam-se em diferentes formalizações do problema e em diferentes técnicas (relaxação Lagrangeana, programação dinâmica com relaxação do espaço dos estados, testes de dominância, etc.) de tratamento dessas formalizações. Em seguida, apresentam-se algumas dessas formalizações, bem como métodos exactos de resolução nelas baseadas.

Formalização 1

O problema de distribuição pode ser considerado uma extensão do problema do caixeiro viajante, como já foi referido. A formalização seguinte, apresentada em Golden et al (1977), baseia-se nessa ideia. Nela são admitidas diferentes capacidades para os diferentes veículos e pode ser generalizada à existência de múltiplos depósitos, Bodin et al (1983).

Sejam:

c_{ij} - distância entre i e j , $i, j = 1, \dots, n$ (o armazém está localizado em 1)

q_i - procura do cliente i , $i = 2, \dots, n$ (pode considerar-se $q_1 = 0$)

Q_k - capacidade do veículo k , $k = 1, \dots, M$

$x_{ijk} = \begin{cases} 1, & \text{se o veículo } k \text{ vai de } i \text{ para } j \text{ directamente} \\ 0, & \text{caso contrário} \end{cases}$

O problema pode formular-se como:

$$\min \sum_{k=1}^M \sum_{i=1}^n \sum_{j=1}^n c_{ij} x_{ijk} \quad (1)$$

$$\text{s. a.} \quad \sum_{i=1}^n \sum_{k=1}^M x_{ijk} \quad j = 1, \dots, n \quad (2)$$

$$\sum_{i=1}^n \sum_{k=1}^M x_{ijk} = 1 \quad i = 1, \dots, n \quad (3)$$

$$\sum_{i=1}^n (x_{ijk} - x_{jik}) = 0 \quad k = 1, \dots, M; j = 1, \dots, n \quad (4)$$

$$\sum_{i=1}^n x_{ijk} \leq 1 \quad k = 1, \dots, M \quad (5)$$

$$\sum_{i=2}^n q_i \left(\sum_{j=1}^M x_{ijk} \right) \leq Q_k \quad k = 1, \dots, M \quad (6)$$

$$x_{ijk} \in \{0, 1\} \quad i, j = 1, \dots, n : k = 1, \dots, M \quad (7)$$

$$X \in S \quad (8)$$

A função objectivo (1) traduz a minimização do comprimento total do trajecto a efectuar pela frota. As condições (2), (3) e (4) garantem que cada cliente é visitado exactamente uma vez e que o veículo que entra numa localização volta a sair dela. As condições (5) asseguram que cada veículo não efectua mais do que uma rota. As condições (6) referem-se às restrições de capacidade nos veículos. As condições (7) definem as variáveis como binárias e as condições (8) são as restrições que evitam a formação de subcircuitos que não passam pelo armazém central e podem ser expressas de diversas formas, Christofides (1979).

Em Magnanti (1981) esta formalização é tratada excluindo as condições (2) (que são desnecessárias na presença das restantes) e relaxando de forma Lagrangeana, Geoffrion (1974), as restrições (3). O problema relaxado pode, então, decompôr-se em M problemas, um para cada veículo. Esta abordagem levanta, contudo, duas questões: a dificuldade de resolução dos problemas relaxados e a possibilidade de na solução dos problemas relaxados ser feita a afectação dum mesmo cliente a vários veículos o que certamente tem más consequências na qualidade do minorante obtido. Neste trabalho não são apresentados resultados computacionais para esta abordagem, nem para as outras variantes sugeridas.

Formalização 2

O problema de distribuição pode também ser relacionado, como já foi referido, com o problema do caixeiro viajante múltiplo. Essa ideia foi explorada por Christofides et al (1981a) num método de pesquisa em árvore em que os minorantes são determinados à custa da resolução de problemas de determinação de árvore com grau k no centro.

Se numa qualquer solução admissível do problema do caixeiro viajante múltiplo, fôr eliminado um conjunto S_0 de y arcos adjacentes ao armazém e um conjunto S_1 de $M - y$ arcos (sendo M o número de caixeiros viajantes) o grafo resultante, se fôr conexo, é uma árvore geradora, Gondran e Minoux (1984), cujo centro (que representa o armazém) tem grau $k = 2M - y$.

Sejam as variáveis:

$$x_{\ell} = \begin{cases} 1, & \text{se o arco } \ell \text{ pertence à árvore de centro com grau } k \\ 0, & \text{caso contrário} \end{cases}$$

$$x_{\ell}^0 = \begin{cases} 1, & \text{se o arco } \ell \text{ pertence ao conjunto } S_0 \\ 0, & \text{caso contrário} \end{cases}$$

$$x_{\ell}^1 = \begin{cases} 1, & \text{se o arco } \ell \text{ pertence ao conjunto } S_1 \\ 0, & \text{caso contrário} \end{cases}$$

O problema do caixeiro viajante múltiplo pode formalizar-se como:

$$\min \sum_{\ell} c_{\ell} (x_{\ell} + x_{\ell}^0 + x_{\ell}^1) \quad (9)$$

$$\text{s.a.} \quad \sum_{\ell \in (V_i, \bar{V}_i)} x_{\ell} \geq 1 \quad \forall V_i \subseteq V \quad (10)$$

$$\sum_{\ell \in A_1} x_{\ell} = 2M - y \quad (11)$$

$$\sum_{\ell} x_{\ell} = n \quad (12)$$

$$\sum_{\ell \in A_0} x_{\ell}^0 = y \quad (13)$$

$$\sum_{\ell \in A_j} (x_{\ell} + x_{\ell}^0 + x_{\ell}^1) = 2 \quad \forall j \neq 1 \quad (14)$$

$$x_{\ell} \in \{0, 1\} \quad \forall_{\ell} \quad (15a)$$

$$x_{\ell}^0 \in \{0, 1\} \quad \forall_{\ell} \quad (15b)$$

$$x_{\ell}^1 \in \{0, 1\} \quad \forall_{\ell} \quad (15c)$$

Com $y \leq M$ e A_i representando os arcos incidentes no vértice i .

A função objectivo (9) representa a minimização do percurso total. As restrições (10) garantem que a solução é conexa. As restrições (11), (12) e (13) impõem o número correcto de arcos na solução. As restrições (14) impõem que cada cliente seja visitado exactamente uma vez e as restrições (15) definem as variáveis como binárias. Relaxando de forma Lagrangeana as restrições (14) obtém-se um minorante para o valor da solução óptima do problema do caixeiro viajante múltiplo (e, portanto, também para o valor da solução óptima do problema de distribuição), o qual é usado numa pesquisa em árvore. As restrições relativas às capacidades dos veículos são apenas tidas em consideração na escolha do valor y e, talvez por isso, os valores obtidos por este processo não são de grande qualidade.

Formalização 3

Uma abordagem diferente do problema de distribuição consiste em considerá-lo como o problema da selecção do subconjunto de rotas, entre todas as rotas possíveis, que minimiza a distância total percorrida.

Seja $R = \{1, 2, \dots, \hat{r}\}$ o conjunto de todas as rotas admissíveis, isto é, que respeitam as restrições de capacidade dos veículos. Para cada cliente i considere-se o conjunto $N_i \subseteq R$ das rotas que o incluem e para cada rota $r \in R$ seja d_r o seu comprimento (ou custo) total.

Sejam as variáveis:

$$y_r = \begin{cases} 1, & \text{se a rota } r \text{ é incluída na solução} \\ 0, & \text{caso contrário} \end{cases}$$

O problema pode ser formalizado como, Balinski e Quandt (1964):

$$\min \sum_{r \in R} d_r y_r \quad (16)$$

$$\text{s.a.} \quad \sum_{r \in N_i} y_r = 1 \quad \forall_i \quad (17)$$

$$y_r \in \{0, 1\} \quad \forall_r \in R \quad (18)$$

A função objectivo (16) traduz a minimização do percurso total. As condições (17) garantem que cada cliente é servido uma vez e as condições (18) definem as variáveis como binárias. Esta formalização corresponde a um problema de partição de um conjunto, podendo portanto o problema de distribuição ser, neste caso, resolvido pelos algoritmos desenvolvidos para o problema da partição de um conjunto, Balas e Padberg (1979). Há, contudo, duas questões de difícil resolução nesta abordagem. Em primeiro lugar, o número total de rotas admissíveis, mesmo para problemas de dimensão moderada, é de tal modo grande que se torna impraticável proceder à sua geração e em seguida resolver o problema de partição de conjunto resultante. Em segundo lugar, a avaliação do comprimento (ou custo) de cada rota exige a resolução de um problema de caixeiro viajante formado pelos clientes inseridos na rota e pelo depósito. Balinski e Quandt optaram por gerar apenas um subconjunto de rotas admissíveis, desenvolvendo desse modo um algoritmo que não garante a obtenção da solução óptima, mas apenas duma solução aproximada.

Christofides et al (1981a) usaram uma formalização idêntica com uma restrição que estabelece a priori o número de rotas na solução. Nela procederam a relaxações de modo a gerar apenas para cada cliente e para cada carga admissível a rota de menor custo, sendo esse custo avaliado por um minorante obtido através do uso da programação dinâmica.

Lucena Filho (1986) com base nesta formalização desenvolveu um conjunto de procedimentos para eliminar à partida uma percentagem substancial de rotas de forma a, trabalhando apenas com um número muito reduzido de rotas, garantir a obtenção da solução óptima. Este trabalho é o que apresenta resultados computacionais para problemas de maior dimensão.

Esta abordagem tem sobre as outras, a vantagem de permitir facilmente a introdução de restrições adicionais ao problema, pois essas restrições diminuem o número de rotas admissíveis.

Formalização 4

Christofides et al (1981b) consideram o problema de distribuição como um problema de programação dinâmica. Admite-se que os clientes e os veículos estão ordenados por ordem decrescente das suas procuras e das suas capacidades, respectivamente. Dado um subconjunto

T de clientes seja $f(k, T)$ o custo mínimo para servir os clientes no conjunto T apenas com os primeiros k veículos, $v(T)$ o valor da solução do problema do caixeiro viajante definido pelo depósito e pelos clientes em T e $q(T)$ a soma das procuras dos clientes em T. A recorrência pode assim definir-se como:

$$f(1, T) = v(T) \quad (19)$$

$$f(k, T) = \min_{S \subset T} \{f(k-1, T-b) + v(S)\} \quad k \geq 2 \quad (20)$$

$$\text{s.a.} \quad q(T) - \sum_{h=1}^{k-1} Q_h \leq q(S) \leq Q_k \quad (21)$$

$$\frac{1}{k-M} q(N-T) \leq q(S) \leq \frac{1}{k} q(T) \quad (22)$$

considerando apenas os subconjuntos T de clientes que satisfazem a:

$$q(N) - \sum_{h=k+1}^M Q_h \leq q(T) \leq \sum_{h=1}^M Q_h \quad (23)$$

As restrições são impostas para evitar o cálculo de valores que não iriam corresponder a rotas admissíveis. Nesta formalização, Christofides et al, procedem à relaxação do espaço dos estados para a obtenção de minorantes a usar num processo de pesquisa em árvore. Os resultados obtidos para os minorantes por este processo são considerados de boa qualidade comparativamente com os obtidos por outras relaxações.

Formalização 5

Uma outra forma de abordar o problema de distribuição é apresentada por Fisher e Jaikumar (1978, 1981). Nesta abordagem o problema é considerado decomponível num problema de afectação generalizada (que determina a afectação dos clientes às rotas) e num conjunto de problemas do caixeiro viajante (que determinam a configuração óptima de cada rota).

Sejam as variáveis:

$$y_{ik} = \begin{cases} 1, & \text{se o cliente } i \text{ é servido pelo veículo } k \\ 0, & \text{caso contrário} \end{cases}$$

$$y_{ijk} = \begin{cases} 1, & \text{se o veículo } k \text{ visita o cliente } j \text{ imediatamente a seguir ao cliente } i \\ 0, & \text{caso contrário} \end{cases}$$

O problema pode formalizar-se como:

$$\min \sum_i \sum_j \sum_k c_{ij} x_{ijk} \quad (24)$$

$$\text{s.a.} \quad \sum_i q_i y_{ik} \leq Q_k \quad k = 1, \dots, M \quad (25)$$

$$\sum_k y_{ik} = \begin{cases} M & i = 1 \\ 1 & i = 2, \dots, n \end{cases} \quad (26)$$

$$y_{ik} \in \{0, 1\} \quad i = 1, \dots, n; \quad k = 1, \dots, M \quad (27)$$

$$\sum_i x_{ijk} = y_{jk} \quad j = 1, \dots, n \quad (28)$$

$$\sum_j x_{ijk} = y_{jk} \quad i = 1, \dots, n \quad (29)$$

$$\sum_{(i,j) \in S \times S} x_{ijk} \leq |S| - 1 \quad S \subset \{2, \dots, n\} \quad (30)$$

$$2 \leq |S| \leq n - 1 \quad (31)$$

$$x_{ijk} \in \{0, 1\} \quad i, j = 1, \dots, n \quad (31)$$

A função objectivo (24) representa a minimização do percurso total. As condições (25), (26) e (27) definem um problema de afectação generalizado, Martello e Toth (1987), que garante a afectação dos clientes às rotas, de acordo com as capacidades dos veículos. As restrições (28), (29), (30) e (31) definem M problemas do caixeiro viajante, um para cada veículo, que garantem a optimização da configuração das rotas. Esta formalização pode ser tomada como base de algoritmos exactos, como sugerido em Fisher e Jaikumar (1978), mas ela é fundamentalmente conhecida como base dum algoritmo heurístico (ver secção seguinte).

Outras Abordagens

Para além dos algoritmos já referidos e que parecem conduzir aos melhores resultados, outras abordagens têm sido sugeridas para o problema de distribuição.

Christofides et al (1979) apresentam um método de pesquisa directa em árvore em que as rotas são geradas sequencialmente por ramificação. Trata-se duma pesquisa em profundidade, em que cada nodo representa uma rota que possa ser efectuada por, pelo menos, um dos veículos da frota e a ramificação é feita escolhendo um cliente ainda não considerado e gerando a lista de todas as rotas admissíveis que incluam esse cliente.

Laporte et al (1985) apresentam uma formalização apenas com dois índices, considerando todos os veículos idênticos, com a qual resolvem problemas simétricos com restrições de capacidade e de distância. A resolução é feita por pesquisa em árvore relaxando as restrições de integralidade e as restrições que impedem a formação de subcircuitos. As primeiras são impostas pela ramificação e as segundas são impostas quando violadas.

Lucena Filho (1986) apresenta uma formalização para o problema básico de distribuição com frota homogênea baseada no problema de fluxo com dois bens e sugere um método de resolução nela baseado. Não apresenta ainda resultados computacionais.

4. Algoritmos Heurísticos

Dada a complexidade computacional dos problemas de distribuição, muito trabalho tem sido também dedicado ao desenvolvimento de algoritmos que permitem obter soluções aproximadamente óptimas com um esforço computacional mais reduzido.

Esses métodos podem ser classificados segundo diversos critérios. Podem ser classificados como construtivos, se se destinam à construção duma solução admissível ou, de tipo melhorativo, se se destinam a melhorar uma solução admissível já conhecida. Os métodos construtivos podem ser classificados de acordo com a estratégia de construção utilizada: sequencial, quando cada rota é construída em separado e só é iniciada a construção de uma nova rota depois de terminada a construção da anterior ou, em paralelo, quando diversas rotas vão sendo construídas simultaneamente. A construção das rotas pode ser feita apenas numa etapa ou em duas. Neste caso, a primeira fase destina-se a agrupar os clientes a incluir em cada rota e, a segunda fase a determinar a sequência pela qual devem ser visitados os clientes incluídos numa mesma rota.

Em seguida, são referidos os principais métodos encontrados na literatura para a resolução aproximada de problemas de distribuição.

(I) Métodos Baseados no Conceito de "Poupança"

O conceito de poupança ("savings" na literatura de língua inglesa) foi inicialmente introduzido por Clarke e Wright (1964) para o problema de distribuição e tem vindo, de então para cá, a ser utilizado em muitos outros problemas. Considere-se que dois clientes i e j , são servidos por dois veículos distintos a partir do armazém central localizado em 0. O custo total envolvido será:

$$c_1 = 2 c_{0i} + 2 c_{0j}$$

Se a capacidade dos veículos permitir que ambos sejam servidos por um só veículo, o custo total envolvido será:

$$c_2 = c_{0i} + c_{0j} + c_{ij}$$

A diferença entre os dois custos

$$s_{ij} = c_{0i} + c_{0j} - c_{ij} \quad (32)$$

é designada por poupança. O algoritmo de Clarke e Wright partindo de uma solução (eventualmente não admissível) em que cada cliente é servido em exclusivo por um veículo, vai procedendo à fusão de rotas correspondentes a poupanças positivas, previamente ordenadas por ordem decrescente, até não ser mais possível reduzir um número de rotas na solução. O algoritmo pode ser implementado, quer na versão sequencial, quer na versão em paralelo.

A fórmula de cálculo da poupança referentes à agregação de duas rotas pode ser modificada. Gaskell (1967) propôs a substituição de (32) por uma das seguintes alternativas

$$\lambda_{ij} = s_{ij} (\bar{c} + |c_{0i} - c_{0j}| - c_{ij}) \quad (33)$$

$$\pi_{ij} = s_{ij} - c_{ij} \quad (34)$$

onde \bar{c} representa a média dos valores c_{0i} e s_{ij} é definido por (32). Yellow (1970) propôs a substituição de (32) por

$$\hat{s}_{ij} = c_{0i} + c_{0j} - \gamma c_{ij} \quad (35)$$

que por manipulação do parâmetro γ permite atribuir maior ou menor importância à distância entre os clientes a agregar numa mesma rota. Em Golden et al (1977) e Paessens (1988) são apresentadas implementações eficientes deste tipo de métodos.

Hart e Shogan (1987) em alternativa à manipulação dos valores da poupança propõem que a selecção das rotas a fundir em cada iteração seja feita de forma aleatória e que o processo seja repetido sucessivas vezes, escolhendo-se no final a melhor solução obtida.

Em Mole e Jameson (1976) é desenvolvido um algoritmo baseado numa generalização deste conceito em que à medida que as rotas vão sendo construídas se vai também considerando a possível alteração da ordem que os clientes nelas ocupam.

(II) Métodos de Duas Fases

Os métodos de duas fases procedem na primeira fase ao agrupamento dos clientes a incluir em cada uma das rotas e na segunda à determinação da configuração de cada rota em função dos clientes que a compõem.

Gillett e Miller (1974) apresentam um algoritmo designado por "sweep algorithm" em que a agregação dos clientes é feita em função das suas localizações representadas pelas respectivas coordenadas polares e a determinação da sequência pela qual são visitados os clientes incluídos numa mesma rota feita com o algoritmo de Lin e Kernighan (1973) para o problema do caixeiro viajante.

Fisher e Jaikumar (1981) apresentam um algoritmo de duas fases baseado na formalização 5, apresentada na secção anterior. Na primeira fase a agregação dos clientes é feita resolvendo um problema de afectação generalizada e na segunda fase é determinada a configuração das rotas, mediante a resolução de problemas do caixeiro viajante com um algoritmo semelhante ao de Millotis (1976, 1978). A resolução do problema de afectação generalizada levanta contudo algumas dificuldades. Trata-se dum problema NP-completo, Martello e Toth (1987), havendo por isso que recorrer para a sua resolução a enumeração em árvore. Em Fisher et al (1986) é apresentado um método de ajustamento de multiplicadores que dá bons resultados na resolução da respectiva relaxação Lagrangeana (neste caso problemas do tipo saco-mochila binário, Martello e Toth (1979)). Para além disso, a função objectivo é neste caso não linear e particularmente difícil de estabelecer (é função das rotas a gerar posteriormente) e é necessário substituí-la por uma aproximação linear.

(III) Métodos de Optimização Incompleta

Os métodos referidos na secção anterior obtêm a solução óptima exacta quando completada a pesquisa em árvore. Eles podem, no entanto, dar origem a métodos aproximados se a pesquisa for restringida.

Balinski e Quandt, como já foi referido, não garantem a obtenção da solução óptima porque restringem a pesquisa a um subconjunto das rotas admissíveis que pode não conter as rotas óptimas.

Christofides et al (1979) sugerem também que o método de pesquisa directa em árvore (ver secção 4) pode ser usado como um método aproximado restringindo o conjunto de rotas considerado e determinando o valor das rotas com métodos heurísticos.

(IV) Outros Métodos

Para além dos métodos referidos nos três pontos anteriores, outros foram também sugeridos. Christofides e Eilon (1969) apresentam um método de troca de arestas semelhante aos métodos r-optimais de Lin (1965) e Lin e Kernighan (1973) para o problema do caixeiro viajante. Os resultados computacionais apresentados são superiores aos obtidos com o método de Clarke e Wright. Christofides et al (1979) apresentam também um método em que numa primeira etapa são escolhidos clientes "semente" para iniciar as rotas e numa segunda etapa se inserem os restantes clientes nas rotas já definidas.

Agradecimento - Este trabalho foi desenvolvido no âmbito do projecto I&D nº 809.86.148/MIC da JNICT

Referências

- BALAS, E. and M. PADBERG (1979), 'Set Partitioning - a Survey', in Combinatorial Optimization, Christofides, N., A. Mingozi, P. Toth and C. Sandi (eds), John Wiley
 BALAS, E. and N. CHRISTOFIDES (1981), 'A Restricted Lagrangean Approach to the Travelling Salesman Problem', Math. Prog., 21, pp 19-46
 BALINSKI, M. and R. QUANDT (1964), 'On an Integer Program for a Delivery Problem', Ops. Res., 12, pp 300-304

- BEASLEY, J. (1984), 'Fixed Routes', *J. Opl. Res. Soc.*, 35, nº 1, pp 49-55
- BODIN, L., B. GOLDEN, A. ASSAD and M. BALL (1983), 'The State of the Art in the Routing and Scheduling of Vehicles and Crews', *Comp. and Ops. Res.*, 10, pp 62-212
- BODIN, L. and T. SEXTON (1986), 'The Multi-Vehicle Subscriber Dial-a-Ride Problem', *Studies in Manag. Sci.*, 22, pp 73-86
- CHRISTOFIDES, N. and S. EILON (1969), 'An Algorithm for the Vehicle-Dispatching Problem', *Opl. Res. Q.*, 20, nº 3, pp 309-318
- CHRISTOFIDES, N. (1979), 'The Travelling Salesman Problem', in *Combinatorial Optimization*, Christofides, N., A. Mingozi, P. Toth and C. Sandi (eds), John Wiley
- CHRISTOFIDES, N., A. MINGOZZI and P. TOTH (1979), 'The Vehicle Routing Problem', in *Combinatorial Optimization*, Christofides, N., A. Mingozi, P. Toth and C. Sandi (eds), John Wiley
- CHRISTOFIDES, N., A. MINGOZZI and P. TOTH (1981a), 'Exact Algorithms for the Vehicle Routing Algorithm Based on Spanning Tree and Shortest Path Relaxations', *Math. Prog.*, 20, pp 255-282
- CHRISTOFIDES, N., A. MINGOZZI and P. TOTH (1981b), 'State-Space Relaxation Procedures for the Computation of Bounds to Routing Problems', *Networks*, 11, pp 145-164
- CLARKE, G. and J. WRIGHT (1964), 'Scheduling of Vehicles from a Central Depot to a Number of Delivery Points', *Ops. Res.*, 12, pp 568-581
- CROWDER, H. and M. PADEBERG (1980), 'Solving Large Scale Symmetric TSP's to Optimality', *Man. Sci.*, 26, pp 495-509
- DESROCHERS, J., Y. DUMAS and F. SOUMIS (1984), 'The Multiple Vehicle Many to Many Routing Problem with Time Windows', *Cahiers du GERARD G-84-13*, Ecole des Hautes etudes de Montreal
- DESROCHERS, J., Y. DUMAS and F. SOUMIS (1986), 'A Dynamic Programming Solution of the Large-Scale Single-Vehicle Dial-a-Ride Problem with Time Windows', *American J. of Math. and Man. Sci.*, 6, nº 3&4, pp 301-325
- DESROSIERS, J.; J.-A. FERLAND, J.-M. ROUSSEAU, G. LAPALME and L. CHAPLEAU (1986), 'TRANSCOL: A Multi-Period School Bus Routing and Scheduling System', *Studies in Man. Sci.*, 22, pp 47-71
- FISHER, M. and R. JAIKUMAR (1978), 'A Decomposition Algorithm for Large Scale Vehicle Routing', WP 78-11-05, Dep. of Decision Sciences, University of Pennsylvania
- FISHER, M. and R. JAIKUMAR (1981), 'A Generalized Assignment Heuristic for the Vehicle Routing', *Networks*, 11, pp 109-124
- FISHER, M., R. JAIKUMAR and W. WASSENHOVE (1986), 'A Multiplier Adjustment Method for the Generalized Assignment Problem', *Man. Sci.*, 32, nº 9, pp 1095-1103
- GASKELL, T. (1967), 'Bases for Vehicle Fleet Scheduling', *Opl. Res. Q.*, 18, pp 281-295
- GAVISH, B. and K. SRIKANTH (1979), 'Mathematical Formulations for the Dial-a-Ride Problem', WP, Graduate School of Management Science, University of Rochester
- GAVISH, B. and S. GRAVES (1979), 'The Travelling Salesman Problem and Related Problems', WP 7906, Graduate School of Management Science, University of Rochester
- GEOFFRION, A. (1974), 'Lagrangian Relaxation for Integer Programming', *Math. Prog. Studies*, 2, pp 82-114
- GILLET, B. and L. MILLER (1974), 'A Heuristic for the Vehicle-Dispatch Problem', *Ops. Res.*, 22, pp 340-349
- GOLDEN, B.; T. MAGNANTI and H. NGUYEN (1977), 'Implementing Vehicle Routing Algorithms', *Networks*, 7, pp 113-148
- GONDRAN, M. and M. MINOUX (1984), *Graphs and Algorithms*, John Wiley
- GRAHAM, D. and H. NUTTE (1986), 'A Comparison for a School Bus Scheduling Problem', *Transp. Res.*, B20, pp 175-182
- HART, J. and A. SHOGAN (1987), 'Semi-Greedy Heuristics: an Empirical Study', *Ops Res Letters*, 6-3, pp 107-114
- LAPORTE, G.; Y. NORBERT and M. DESROCHERS (1985), 'Optimal Routing under Capacity and Distance Constraints', *Ops. Res.*, 33, pp 1050-1073
- LAPORTE, G. and Y. NORBERT (1987), 'Exact Algorithms for the Vehicle Routing Problem', *Annals of Disc. Math.*, 31, pp 147-184
- LENSTRA, J. and A. RINNOOY KAN (1981), 'Complexity of Vehicle Routing and Scheduling Problems', *Networks*, pp 221-227
- LIN, S. (1965), 'Computer Solutions of the Travelling Salesman Problem', *The Bell Syst. Techn. J.*, pp 2245-2269

- LIN, S. and B. KERNIGHAN (1973), 'An Effective Heuristic Algorithm for the Traveling Salesman Problem', *Ops. Res.*, 21, pp 498-516
- LUCENA FILHO, A. (1986), 'Exact Solution Approaches for the Vehicle Routing Problem', Ph. D. Thesis, University of London
- MAGNANTI, T. (1981), 'Combinatorial Optimization and Vehicle Fleet Planning Perspectives and Prospects', *Networks*, 11, pp 179-214
- MARTELLO, S. and P. TOTH (1979), 'The 0-1 Knapsack Problem', In *Combinatorial Optimization*, Christofides, N., A. Mingozzi, P. Toth and C. Sandi (eds), Wiley
- MARTELLO, S. and P. TOTH (1987), 'Linear Assignment Problems', *Annals of Disc. Math.*, 31, pp 259-282
- MILIOTIS, P. (1976), 'Integer Programming Approaches to the Travelling Salesman Problem', *Math Prog.*, 10, pp 367-378
- MILIOTIS, P. (1978), 'Using Cutting Planes to Solve the Symmetric Travelling Salesman Problem', *Math. Prog.*, 15, pp 177-188
- MOLE, R. and S. JAMESON (1976), 'A Sequential Route-Building Algorithm Employing a Generalised Savings Criterion', *Opl. Res. Q.*, 27, pp 503-511
- PADBERG, M. and G. RINALDI (1986), 'Optimization of a 532-City Symmetric Travelling Salesman Problem', *Journee du 20-ème Anniversaire du Groupe Combinatoire, AFCET*, 3,4,5 Dec.
- PAESSENS; H. (1988), 'The Savings Algorithm for the Vehicle Routing Problem' , *EJOR*, 34, pp 336-344
- SOLOMON, M. and J. DESROSIERS (1987), 'Time Window Constrained Routing and Scheduling Problems: a Survey', *Cahiers du GERARD G-87-09, Groupe d'Etudes et de Recherche en Analyse des Décisions, Montreal*
- YELLOW; P. (1970), 'A Computational Modification to the Savings Method of Vehicle Scheduling', *Opl. Res. Q.*, 21, pp 281-283

CUTTING PLANES FROM CONDITIONAL BOUNDS FOR GENERALIZED SET COVERING PROBLEMS

Margarida V. Pato
Dept. de Matemática, ISEG - UTL
Rua Miguel Lupi, 20
1200 LISBOA

José M. P. Paixão
DEIOC - CEAUL (INIC)
Av. 24 de Julho, 134-5º
1300 LISBOA

Abstract. This paper reports on the development of special cutting planes for the generalized set covering problem, GSCP, which is a covering problem where the variables and the right-hand sides are allowed to have any positive integer value. Those inequalities are, actually, a generalization of the cutting planes derived from conditional bounds and originally presented by Balas (1980), for the set covering problem. More recently, Hall & Hochbaum (1985) have extended those results for the multicovering problem. The generalized inequalities that we derive for the GSCP are proved to be of the covering type and, hence, keeping the structure of the problem constraints.

Key Words: cutting planes, conditional bounds, disjunctive inequalities, generalized set covering problem

1. Introduction

In this paper, we deal with the generalized set covering problem (GSCP), which can be stated as the following mathematical program:

$$\begin{aligned}
 \text{(GSCP)} \quad & \min \sum_{j \in N} c_j x_j \\
 \text{s. to} \quad & \sum_{j \in N} a_{ij} x_j \geq b_i \quad (i \in M) \quad (1) \\
 & 0 \leq x_j \leq h_j \text{ and integer} \quad (j \in N) \quad (2)
 \end{aligned}$$

where M and N are the index sets of, respectively, the rows and columns for the problem. The b_i ($i \in M$) and h_j ($j \in N$) are positive integer values. Also, one has $a_{ij} \in (0, 1)$ ($i \in M, j \in N$) and, in order to avoid a trivial resolution, we assume that $\sum_{j \in N} a_{ij} > 0$. A vector x verifying the constraints (1) and (2) is called a cover. Finally, and for the sake of simplicity of the notation, we define the following two sets: $M_j = \{i \in M : a_{ij} = 1\}$ and $N_i = \{j \in N : a_{ij} = 1\}$. Then, for instance, the covering constraints can be stated as $\sum_{j \in N_i} x_j \geq b_i$ ($i \in M$).

The GSCP has been widely used for several real life situations and a survey for that can be found in Pato (1989). Most of those applications are related to personnel scheduling, particularly the determining of the schedules for drivers in mass transit bus companies [Blais & Rousseau (1982); Bodin, Rosenfield & Kydes (1981); Mitra & Welsh (1981); Paixão et al (1986); Shepardson & Marsten (1980); Yihua (1985)].

As suggested by the designation, the GSCP includes the well known set covering problem (SCP), where all the b_i ($i \in M$) and h_j ($j \in N$) are equal to the unity. It also includes the so called multicovering problem (MCP), where $h_j = 1$ ($j \in N$) but the right-hand side values b_i ($i \in M$) can be any positive integers. From that, it follows straightforwardly that the GSCP is an NP-hard problem.

Hence, heuristics and LP-based techniques have been the approaches most used for the GSCP. In particular, cutting plane methods have been applied as a way of dealing with instances where a null linear gap occurs or a high number of alternative solutions exists [Geoffrion & Marsten (1972); Wolfe (1984)]. That, actually, is the case of GSCPs related to crew scheduling problems [Pato (1989)], for which lagrangean relaxation and tree-search procedures have been applied too [Paixão & Pato (1989); Shepardson & Marsten (1980)].

In the present paper, we extend to the GSCP a class of cutting planes that have been described by Balas (1980) and Balas & Ho (1980) for the SCP. Those cutting planes are derived from conditional bounds following a general disjunctive approach [Balas (1975), (1979)], and have been extended by Hall & Hochbaum (1985) for the MCP.

The paper is organized as follows. Next, in this section, we introduce the idea of cutting planes from conditional bounds through an example. Some formal notation is stated too. Section 2 is devoted to a single result leading to the obtaining of a valid inequality which is strengthened in the following section. Then, in section 4, the example considered in the introduction is used for the purpose of illustrating the results of the previous sections. Finally, some remarks and conclusions are presented in section 5.

Before giving an example for introducing the idea of conditional bounds and the related cutting planes, let us state some notation. We denote by $\overline{\text{GSCP}}$ the continuous version of the GSCP, and the corresponding dual linear problem, $\overline{\text{DGSCP}}$, has the following formulation:

$$\begin{aligned}
 (\overline{\text{DGSCP}}) \quad & \max \quad \sum_{i \in M} b_i u_i - \sum_{j \in N} h_j v_j \\
 \text{s. to} \quad & \sum_{i \in M_j} u_i - v_j \leq c_j \quad (j \in N) \quad (3) \\
 & v_j \geq 0 \quad (j \in N) \quad (4) \\
 & u_i \geq 0 \quad (i \in M) \quad (5)
 \end{aligned}$$

The constraints (3) can be rewritten as $r_j = c_j - \sum_{i \in M_j} u_i + v_j \geq 0$ ($j \in N$), with r_j being designated as the reduced cost for the variable index j . A vector $[u \ v]$ satisfying constraints (3)-(5) is said to be a dual feasible solution.

Now, let us consider the following instance for the GSCP:

$$\begin{aligned}
 \min \quad & 2x_1 + 2x_2 + x_3 + 1.5x_4 + 2x_5 \\
 \text{s. to} \quad & x_2 + x_3 + x_5 \geq 4 \\
 & x_1 + x_4 + x_5 \geq 2 \\
 & x_2 + x_4 + x_5 \geq 6 \\
 & x_1 \leq 1, x_2 \leq 2, x_3, x_4, x_5 \leq 4 \\
 & x_1, \dots, x_5 \geq 0 \text{ and integers.}
 \end{aligned}$$

Also, consider a cover given by $\tilde{x} = [0 \ 2 \ 2 \ 4 \ 0]$ and the dual feasible solution defined by $\tilde{u} = [0 \ 0 \ 1.5]$ and $\tilde{v} = [0 \ 0 \ 0 \ 0 \ 0]$. Hence, the optimum is a value between 9 and 12.

Suppose that the constraint $x_3 \geq 3$ is added to the covering problem. Then, for this enlarged instance, one may consider the previous dual feasible solution with an additional variable $\tilde{u}_4 = 1$. The optimal value for the enlarged instance is bounded from below by 12. This leads to the

conclusion that the constraint $x_3 < 3$ has to be satisfied by any feasible solution for the original problem with an objective value less than 12. Combining this with the first constraint of the covering problem, $x_2 + x_3 + x_5 \geq 4$, one may conclude that $x_2 + x_5 \geq 2$ is a valid inequality for any cover with a better value than the current one. Note, that this new constraint is of the covering type.

2. Weak Inequality

In this section, we derive a first valid inequality for the GSCP following the approach proposed by Balas & Ho (1980) for the set covering problem. For that, a cover and a dual feasible solution must be available. Consequently, an upper bound z_u for the optimal value of the GSCP is known. Then, as suggested in the example, additional constraints are determined in such a way that the corresponding dual linear problem has an optimal value greater than or equal to z_u . Therefore, solutions strictly better than z_u must violate at least one of those additional constraints. Finally, the combining of this last disjunction to some of the covering constraints for the original GSCP leads to a valid inequality for the feasible solutions with a better value than the current one.

The required pair of feasible solutions, \tilde{x} for the GSCP and $[\tilde{u} \ \tilde{v}]$ for the dual linear problem (DGSCP), must verify the following conditions:

$$\left(\sum_{j \in N_i} \tilde{x}_j - b_i \right) \tilde{u}_i = 0 \quad (i \in M) \quad (6)$$

$$\tilde{v}_j = \max \left\{ 0, -c_j + \sum_{j \in M_i} \tilde{u}_i \right\} \quad (j \in N) \quad (7)$$

$$(\tilde{x}_j - h_j) \tilde{v}_j = 0 \quad (j \in N) \quad (8)$$

We denote by \tilde{z}_u and \tilde{z}_ℓ , respectively, the upper bound on the optimum provided by \tilde{x} and the lower bound given by the associated dual feasible solution $[\tilde{u} \ \tilde{v}]$. Those solutions can be easily obtained through the using of primal-dual greedy heuristics presented for the GSCP by Paixão & Pato (1987).

Now, let us define the set $\tilde{S} = \{j \in N : \tilde{x}_j > 0 \text{ and } \tilde{r}_j > 0\}$ with \tilde{r}_j being the reduced cost produced by $[\tilde{u} \ \tilde{v}]$ for the j th column. And, let $S = \{j(1), \dots, j(p)\}$ be a subset of \tilde{S} and integers $\delta_{j(k)} \geq \tilde{x}_{j(k)}$ ($j(k) \in S$) such that:

$$\sum_{k=1, \dots, p} \tilde{r}_{j(k)} \delta_{j(k)} \geq z_u - \tilde{z}_\ell \quad (9)$$

where z_u is the best known upper bound on the optimum for GSCP.

Note that such condition holds for the case where $S = \tilde{S}$, $z_u = \tilde{z}_u$ and $\delta_{j(k)} = \tilde{x}_{j(k)}$.

Theorem 1. Let \tilde{x} and $[\tilde{u} \ \tilde{v}]$ be feasible solutions, respectively, for GSCP and DGSCP verifying the conditions (6)-(8). Also, let $\delta_{j(k)}$ ($k=1, \dots, p$) be integer values for which (9) holds with z_u , \tilde{z}_ℓ and \tilde{r} defined as above.

Then each feasible solution, x , whose value is less than z_u must satisfy

$$\sum_{j \in W} x_j \geq d \quad (10)$$

with

$$d = \min_{k=1, \dots, p} (b_{j(k)} - \delta_{j(k)}) + 1 \quad (11)$$

$$W = \bigcup_{k=1, \dots, p} (N_{i(k)} - Q_k) \quad (12)$$

where, $i(k)$ is any index in M and $Q_k \subseteq N_{i(k)}$ ($k=1, \dots, p$) verify

$$\sum_{k: j \in Q_k} \tilde{r}_{j(k)} \leq \tilde{r}_j \quad (j \in N) \quad (13)$$

Proof. Consider $S = \{j(1), \dots, j(p)\}$ and a vector δ , both defined according to our hypothesis. Let GSCP_A denote the GSCP enlarged with the p additional constraints:

$$\sum_{j \in Q_k} x_j \geq \delta_{j(k)} \quad (k=1, \dots, p) \quad (14)$$

The dual linear of this enlarged problem has a feasible solution, given by $[\bar{u} \ \bar{v}]$ plus p variables associated to new constraints (14) and, respectively, equal to the reduced costs, $\bar{r}_j(1), \dots, \bar{r}_j(p)$. The definition of the sets Q_k ($k=1, \dots, p$) guarantees the dual feasibility for this 'enlarged' solution.

The value of this dual feasible solution for the enlarged problem is given by

$$\bar{z}_d + \sum_{k=1, \dots, p} \bar{r}_j(k) \delta_{j(k)}$$

which, according to (9), is greater than or equal to z_u , the upper bound for the optimal value.

From the weak duality theorem applied to the GSCP_A, one may conclude that the enlarged problem has no feasible solution better than z_u . Then any feasible solution of GSCP better than z_u must violate at least one of the additional constraints (14). In other words, such cover must satisfy the following disjunction

$$\bigvee_{k=1, \dots, p} \left(\sum_{j \in Q_k} x_j < \delta_{j(k)} \right),$$

which implies

$$\bigvee_{k=1, \dots, p} \left(\sum_{j \in N_{i(k)} - Q_k} x_j \geq b_{i(k)} - \delta_{j(k)} + 1 \right) \quad (15)$$

Now, if d and W are defined by (11)-(13), one has that $\sum_{j \in W} x_j \geq d$ is an inequality satisfied

by any solution with a value better than z_u . ◇

The existence of the sets Q_k ($k = 1, \dots, p$) is guaranteed and, next, we present a procedure (denoted by QKAPA), which produces such a family of sets. There, we assume that Criterium j_r and Criterium i_s have been established, respectively, for selecting a row and a column in each iteration.

Using different criteria for determining $j(k)$ and $i(k)$ in the procedure QKAPA, several cuts can be generated. Naturally, if we intend to obtain a strong cut, it is reasonable to keep the cardinality of W as low as possible and, simultaneously, find a large value for d , the right-hand side of the inequality.

Hence, in order to reduce $|W|$ one should try to add the least possible number of additional constraints. This can be pursued through including the columns in S by decreasing order of the reduced costs (Criterium j_1). Of course, if $z_u = \bar{z}_u$ and $\delta_{j(k)} = \bar{x}_{j(k)}$ then p is equal to $|S|$ and any sorting of the reduced costs is irrelevant.

Another way of reducing $|W|$ consists of including as many variables as possible in the additional constraints, while respecting dual feasibility. The reduction of $|W|$ may also be achieved by choosing, accordingly to Criterium i_2 , $i(k)$ as the row with minimum cardinality among the ones covered by column $j(k)$.

The determination of a large value d could be attained by taking $\delta_{j(k)} = \bar{x}_{j(k)}$ and by choosing row $i(k)$ in $M_{j(k)}$ such that $b_{i(k)}$ is the largest one (Criterium i_1).

Therefore, the selection rules for column $j(k)$ and row $i(k)$ may be relevant. Unfortunately these rules are not enough to ensure the elimination of any feasible solution for GSCP, which is the main objective when deducing efficient cutting planes. In the next section we present and discuss a way of achieving that.

```

procedure QKAPA ( input: GSCP,  $\bar{r}$ ,  $S$  ; output:  $Q_k$  ( $k=1, \dots, p$ ) )
  * initialization *
   $r^A \leftarrow \bar{r}$  ;  $p \leftarrow |S|$ 
  * iterations *
  for  $k=1, \dots, p$  do
    choose  $j(k)$  applying Criterium  $j_r$ 
     $Q_k \leftarrow \{j(k)\}$ 
    choose  $i(k) \in M_{j(k)}$  applying Criterium  $i_s$ 
    for  $j \in N_{i(k)} - S$  do
      if  $r_j^A \geq \bar{r}_{j(k)}$  then ( $Q_k \leftarrow Q_k \cup \{j\}$ ) endif
    enddo
     $r_j^A \leftarrow r_j^A - \bar{r}_{j(k)}$  ( $j \in Q_k$ )
  enddo
end QKAPA

```

3. Strong Inequality

In this section, we define a cutting plane from conditional bounds that eliminates a previously known nonredundant cover (one whose components cannot be reduced without producing unfeasibilities). Also, it will be proved that such cut defines a facet of a GSCP polyhedron.

First, let us define the following sets of variables:

$$S^h = \{j \in N : \tilde{x}_j = h_j\} \tag{16}$$

$$T = \{i \in M : \sum_{j \in N_i} \tilde{x}_j = b_i\} \tag{17}$$

Associated to a given variable $j(k)$ we define a row:

$$i(k) \in T \cap M_{j(k)} \text{ such that } \tilde{x}_j = 0 \text{ for } j \in N_{i(k)} - S^h - \{j(k)\} \tag{18}$$

In order to build up the sets Q_k ($k=1, \dots, p$), we shall now consider this last condition in the row selection criterium of the procedure QKAPA.

Theorem 2. Let \tilde{x} be a nonredundant cover for the GSCP, the sets S^h , Q_k ($k=1, \dots, p$) and indices $i(k)$ ($k=1, \dots, p$) be defined as above. Also, let $z_u, \tilde{z}_u, p, \{\tilde{u} \ \tilde{v}\}, \tilde{z}_d, \tilde{r}$ stand as in Theorem 1. Then every feasible solution to the GSCP whose value is less than z_u verifies the following condition violated by \tilde{x} :

$$\sum_{j \in \bar{W}} x_j \geq 1 \tag{19}$$

where

$$\bar{W} = \bigcup_{k=1, \dots, p} (N_{i(k)} - Q_k - S^h) \tag{20}$$

Proof. Taking into account that $S = \tilde{S}$ and $\delta = \tilde{x}$, the assumptions of Theorem 1 are verified for any $z_u \leq \tilde{z}_u$. Consider that the procedure QKAPA determines Q_k ($k=1, \dots, p$) with the further condition in selecting row $i(k)$, given by (18).

As was seen in the proof of the first theorem, the improved solutions relatively to z_u violate at least one of those additional constraints. That is, they satisfy the condition (15).

Now, let us rearrange the condition (15) in the following form:

$$\bigvee_{k=1, \dots, p} \left(\sum_{j \in N_{i(k)} - Q_k - S^h} x_j \geq b_{i(k)} - \tilde{x}_{j(k)} + 1 - \sum_{j \in (N_{i(k)} - Q_k) \cap S^h} h_j \right) \tag{21}$$

Note that, from the definition, $\tilde{x}_j = h_j$ for all $j \in (N_{i(k)} - Q_k) \cap S^h$.

Since $i(k) \in T$ and (18) is verified for all $k=1, \dots, p$, the above condition (21) is equivalent to

$$\bigvee_{k=1, \dots, p} \left(\sum_{j \in N_{i(k)} - Q_k - S^h} x_j \geq 1 \right)$$

Now from the integrality of the GSCP variables, one has the following inequality:

$$\sum_{j \in \bigcup_{k=1, \dots, p} (N_{i(k)} - Q_k - S^h)} x_j \geq 1 \tag{22}$$

which is valid for any feasible solution strictly better than z_u .

Since $\tilde{x}_j = 0$ for $j \in \bigcup_{k=1, \dots, p} (N_{i(k)} - Q_k - S^h)$, it comes straightforwardly that \tilde{x} is eliminated

by the deduced cut ◊

Note that the last result states that the cut (19) is a valid inequality for the set of all $|N|$ -dimensional vectors satisfying constraints (1) - (2) amended with $\sum_{j \in N} c_j x_j < \tilde{z}_u$.

Now, one may generalize, for the GSCP, the known facet defining property of the similar cut derived for the SCP [Balas (1980)]. Before stating that, let P be the convex hull of all integer nonnegative $|N|$ -dimensional vectors satisfying (1) and (19). That is,

$$P = \text{conv} \left\{ x \in \mathbb{R}^{|N|} : \sum_{i \in M} x_i \geq b_i \ (i \in M), \sum_{j \in \bar{W}} x_j \geq 1 \text{ and } x_j \text{ nonnegative integer } (j \in N) \right\} \tag{23}$$

Theorem 3. The cutting plane (19), constructed according to the Theorem 2, defines a facet of the polyhedron P given by (23).

Proof. First let us prove that the index set \bar{W} , defined through the hypothesis of Theorem 2, satisfies the following:

$$\nexists k \in M : N_k \subseteq \bar{W} \tag{24}$$

Consider \tilde{x} , the cover for the GSCP defined in that theorem, and the set of its positive variables, $N_{\tilde{x}} = \{ j \in N : \tilde{x}_j > 0 \}$. Thus, one has $N_k \cap N_{\tilde{x}} \neq \Phi$ for all $k \in M$.

But \tilde{x} violates the inequality (19) and, therefore $N_{\tilde{x}} \cap \bar{W} = \Phi$. Then, the statement (24) is easily seen to be true. This will be used later on.

We already know that (19) is a valid inequality for the polyhedron P. Now, one aims to build up $|N|$ linearly independent $|N|$ -dimensional nonnegative integer vectors satisfying (1) and verifying (19) as an equality.

Let us assume without loss of generality that \bar{W} is the set of the first p indexes in N . Consider the $|N|$ vectors given by the rows of the following matrix:

$$X = \begin{array}{c} \begin{array}{cccccccc} \leftarrow & \dots & p & \dots & \rightarrow & \leftarrow & \dots & |N| - p & \dots & \rightarrow \\ 1 & 0 & 0 & \dots & 0 & 0 & \ell & \ell & \ell & \dots & \ell & \ell \\ 0 & 1 & 0 & \dots & 0 & 0 & \ell & \ell & \ell & \dots & \ell & \ell \\ 0 & 0 & 1 & \dots & 0 & 0 & \ell & \ell & \ell & \dots & \ell & \ell \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ 0 & 0 & 0 & \dots & 1 & 0 & \ell & \ell & \ell & \dots & \ell & \ell \\ 0 & 0 & 0 & \dots & 0 & 1 & \ell & \ell & \ell & \dots & \ell & \ell \\ 1 & 0 & 0 & \dots & 0 & 0 & 2\ell & \ell & \ell & \dots & \ell & \ell \\ 1 & 0 & 0 & \dots & 0 & 0 & \ell & 2\ell & \ell & \dots & \ell & \ell \\ 1 & 0 & 0 & \dots & 0 & 0 & \ell & \ell & 2\ell & \dots & \ell & \ell \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ 1 & 0 & 0 & \dots & 0 & 0 & \ell & \ell & \ell & \dots & 2\ell & \ell \\ 1 & 0 & 0 & \dots & 0 & 0 & \ell & \ell & \ell & \dots & \ell & 2\ell \end{array} & \begin{array}{c} \uparrow \\ \vdots \\ \downarrow \\ \uparrow \\ \vdots \\ \downarrow \\ \uparrow \\ \vdots \\ \downarrow \end{array} \end{array}$$

where $\ell = \max_{i \in M} b_i$.

Taking into account that (24) holds true for our case, each one of these vectors verifies the constraints (1). Moreover, it has nonnegative integer components and it satisfies strictly inequality (19), once there is only one 1 in its first p components.

It now remains to show that X is a nonsingular matrix.

Let us define matrix Z by:

$$\begin{array}{c} \begin{array}{cccccccc} \leftarrow & \dots & p & \dots & \rightarrow & \leftarrow & \dots & |N| - p & \dots & \rightarrow \\ 1+|N|-p & 0 & 0 & \dots & 0 & 0 & -1 & -1 & -1 & \dots & -1 & -1 \\ |N|-p & 1 & 0 & \dots & 0 & 0 & -1 & -1 & -1 & \dots & -1 & -1 \\ |N|-p & 0 & 1 & \dots & 0 & 0 & -1 & -1 & -1 & \dots & -1 & -1 \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ |N|-p & 0 & 0 & \dots & 1 & 0 & -1 & -1 & -1 & \dots & -1 & -1 \\ |N|-p & 0 & 0 & \dots & 0 & 1 & -1 & -1 & -1 & \dots & -1 & -1 \\ -1/\ell & 0 & 0 & \dots & 0 & 0 & 1/\ell & 0 & 0 & \dots & 0 & 0 \\ -1/\ell & 0 & 0 & \dots & 0 & 0 & 0 & 1/\ell & 0 & \dots & 0 & 0 \\ -1/\ell & 0 & 0 & \dots & 0 & 0 & 0 & 0 & 1/\ell & \dots & 0 & 0 \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ \cdot & \cdot & \cdot & \dots & \cdot & \cdot & \cdot & \cdot & \cdot & \dots & \cdot & \cdot \\ -1/\ell & 0 & 0 & \dots & 0 & 0 & 0 & 0 & 0 & \dots & 1/\ell & 0 \\ -1/\ell & 0 & 0 & \dots & 0 & 0 & 0 & 0 & 0 & \dots & 0 & 1/\ell \end{array} & \begin{array}{c} \uparrow \\ \vdots \\ \downarrow \\ \uparrow \\ \vdots \\ \downarrow \\ \uparrow \\ \vdots \\ \downarrow \end{array} \end{array}$$

and see that $XZ = I_{|N|}$, where $I_{|N|}$ stands for the $|N| \times |N|$ identity matrix.

Hence, the $|N|$ vectors given above are linearly independent and the result comes straightforwardly from the definition of facet of a $|N|$ -dimensional polyhedron. \diamond

In the next page we describe the procedure CUT which, when successful, produces a cutting plane under the conditions stated above.

However, the procedure CUT may fail in obtaining the 'strong' cut that one aims for. This is due to the fact that, for most of the instances, the hypothesis of Theorem 2, namely the one expressed by the condition (18), can hardly be verified. This is completely different from the case of the SCP and MCP where the variable upper bounds are all equal to the unity. For those cases, the condition (18) is easily satisfied and the cuts from conditional bounds can be more effective [Balas & Ho (1980); Fernandez (1985); Hall & Hochbaum (1985)].

```

procedure CUT ( input: GSCP,  $z_u, \bar{x}, \bar{z}_u, [\bar{u} \bar{v}], \bar{z}_l$  and  $\bar{r}$  satisfying the hypothesis of Theorem 2 ;
output: should the answer be true, cut  $\sum_{i \in \bar{W}} x_j \geq 1$  )

* initialization *
D  $\leftarrow \bar{z}_l$  ; k  $\leftarrow 1$  ;  $\bar{W} \leftarrow \Phi$  ;  $\bar{S} \leftarrow \{j \in N : \bar{x}_j \bar{r}_j > 0\}$  ;  $S^0 \leftarrow \{j \in N : \bar{x}_j = 0\}$ 
 $S^h \leftarrow \{j \in N : \bar{x}_j = h_j\}$  ; T  $\leftarrow \{i \in M : \sum_{i \in N_i} \bar{x}_j = b_i\}$  ;  $r_j^A \leftarrow \bar{r}_j$  (j  $\in N$ ) ; answer  $\leftarrow$  true

* iterations *
label 1: j(k)  $\leftarrow$  arg max  $\bar{r}_j$  ; * Criterium j1 *
j  $\in \bar{S}$ 

 $\bar{O}_k \leftarrow \Phi$  ;  $\bar{T} \leftarrow \{i \in T \cap M_{j(k)} : (N_i - S^h - \{j(k)\}) \cap S^0 = (N_i - S^h - \{j(k)\})\}$ 
if  $\bar{T} \neq \Phi$ 
then choose i(k)  $\in \bar{T}$  with Criterium  $i_s$  ; * select row index *
else answer  $\leftarrow$  false ; stop ; * cutting plane not determined *
endif

for j  $\in (N_{i(k)} - \bar{S} - \bar{W}) \cup \{j(k)\}$  do
if  $r_j^A < \bar{r}_{j(k)}$ 
then ( if  $\bar{x}_j = 0$  then  $\bar{O}_k \leftarrow \bar{O}_k \cup \{j\}$  endif ) * variables for the cut *
else  $r_j^A \leftarrow r_j^A - \bar{r}_{j(k)}$  * variables for the new additional constraint *
endif
enddo

D  $\leftarrow D + \bar{r}_{j(k)} \bar{x}_{j(k)}$  ;  $\bar{W} \leftarrow \bar{W} \cup \bar{O}_k$  ; * verify validity of the cut *
if D <  $z_u$  then (  $\bar{S} \leftarrow \bar{S} - \{j(k)\}$  ; k  $\leftarrow k + 1$  ; goto label 1 ) endif

if  $\bar{W} = \Phi$  then answer  $\leftarrow$  false endif ; * optimal solution for GSCP found *
end CUT
    
```

A further difficulty, in practical terms, comes from the fact that a valid inequality obtained from Theorem 2 eliminates the current feasible solution, but does not necessarily cut off other solutions with the same value. And this is very frequent for the GSCP.

Thus the procedure without improvement turns out to be very slow. In fact, our personal experience in using this approach for real life GSCPs proved to be very poor. The procedure CUT was tried out for 34 instances related to a real life scheduling application, all of them with 36 rows, 100 or 865 columns and high density - more than 50% of 1s in the covering constraint matrix [Paixão & Pato (1989); Pato (1989)]. The primal and dual solutions needed for generating the cuts were found by means of a combined primal-dual greedy and improving heuristic procedure [see Paixão & Pato (1987)]. However, we were not successful with any one of the above test problems.

4. Example

To illustrate the results of the previous sections we consider the example given in section 1 assuming that the same pair of feasible solutions is available and also that $z_u = \bar{z}_u = 12$. In this case, one has that $\bar{r} = [2 \ 0.5 \ 1 \ 0 \ 0.5]$ and $\bar{z}_l = 9$. Since $\bar{r}_3 \delta_3 = 1 \times 3$ and $\bar{z}_u - \bar{z}_l = 12 - 9$, we may easily verify that the hypothesis of Theorem 1 is true, with $S = \{3\}$ and $p = 1$.

Let us take $Q_1 = 3$, $i(1)=1$ and, then, calculate $W = N_{i(1)} - Q_1 = N_1 - Q_1 = \{2,5\}$ and $d = b_{i(1)} - \delta_{j(1)} + 1 = b_1 - \delta_3 + 1 = 4 - 3 + 1 = 2$.

The previously constructed valid inequality, $x_2 + x_5 \geq 2$, may now be derived in accordance with the first theorem.

Bearing in mind the same cover and the same dual feasible solution, along with the fact that condition (9) is fulfilled with $S = \bar{S}$ and $\delta_{j(k)} = \bar{x}_{j(k)}$, we can see that $S = \{2,3\}$ and so $p=2$, in the hypothesis of Theorem 2.

Firstly, the column $j(1) = 2$ is chosen and any row selection criterium picks up $i(1) = 3$, because it is the only one from set $M_2 = \{1,3\}$ satisfying condition (18). Thus $Q_1 = \{2,5\}$ and the reduced costs are updated: $r^A = [2 \ 0 \ 1 \ 0 \ 0]$. At last $j(2) = 3$ and now $i(2) = 1$, the only possible choice, leading to $Q_2 = \{3\}$. The reduced costs for the enlarged problem become $r^A = [2 \ 0 \ 0 \ 0 \ 0]$.

As the set of variables equal to their upper bounds is $S^h = \{2,4\}$, the variables with positive coefficient in the cut belong to the set $\bar{W} = \{N_3 - Q_1 - S^h\} \cup \{N_1 - Q_2 - S^h\} = \{5\}$.

Thus, from Theorem 2, the valid cutting plane is $x_5 \geq 1$, which is clearly violated by the current cover and is verified by every strictly less-than-12 cover. As may be observed, the three optimal alternative solutions with value 11 ($x_4 = 2, x_5 = 4; x_2 = 1, x_4 = 2, x_5 = 3; x_2 = x_4 = x_5 = 2$) satisfy this new valid inequality.

5. Remarks

In this paper, we have characterized two families of valid inequalities for the GSCP. These results correspond to a generalization of Balas's cuts for the SCP. Those cuts keep the covering structure.

The first type, studied throughout section 2, can be defined from conditional bounds for every GSCP, but does not necessarily eliminate one single feasible solution.

The other kind of cutting plane, also derived from conditional bounds (Theorem 2), eliminates the current cover. However, very tight conditions are required and, thus, one gets a harder process to generate the cuts for genuine GSCP instances. In fact, our computational experience attested this difficulty.

The study of a less restrictive hypothesis in Theorem 2 could be a promising field of research, and we hope that an algorithm combining these cuts with several heuristic or other bounding tools will enable us to efficiently tackle the generalized set covering problem.

Acknowledgments

The authors are grateful to Profs. Jaume Barceló and Elena Fernandez of Facultat d'Informàtica of Universitat Politècnica de Catalunya for their suggestions at the beginning of this study.

References

- E. Balas (1975) "Disjunctive programming: Cutting planes from logical conditions", in *Nonlinear Programming 2*, O. L. Mangasarian et al.(eds.), Academic Press, pp. 279-312.
- E. Balas (1979), "Disjunctive programming", *Annals of Discrete Mathematics*, vol. 5, pp. 3-51.
- E. Balas (1980), "Cutting planes from conditional bounds: A new approach to set covering", *Mathematical Programming*, vol. 12, pp. 19-36.
- E. Balas and A. Ho (1980), "Set covering algorithms using cutting planes, heuristics and subgradient optimization: A computational study", *Mathematical Programming*, vol. 12, pp. 37-60.
- J.-Y. Blais and J.-M. Rousseau (1982), "HASTUS: A model for the economic evaluation of drivers' collective agreements in transit companies", *INFOR*, vol. 20, pp. 3-15.
- L. D. Bodin, D. B. Rosenfield and A. S. Kydes (1981), "Scheduling and estimation techniques for transportation planning", *Computers & Operations Research*, vol. 8, pp. 25-38.
- E. Fernandez (1985), *Experiencias Computacionales con Algoritmos de Planos Secantes para Problemas de Set Covering*, Tesina, Facultat de Matemàtiques, Universidad de Valencia, Valencia.

- A. M. Geoffrion and R. E. Marsten (1972), "Integer programming algorithms: A framework and state-of-the-art survey", *Management Science*, vol. 18, pp. 465-491.
- N. G. Hall and D. S. Hochbaum (1985), "The multicovering problem: The use of heuristics, cutting planes and subgradient optimization for a class of integer programs", Working Paper WPS 85-73, College of Administrative Science, The Ohio State University, Columbus.
- G. Mitra and A. Welsh (1981), "A computer based crew scheduling system using a mathematical programming approach", in *Computer Scheduling of Public Transport: Urban Passenger Vehicle and Crew Scheduling*, A. Wren (ed.), North-Holland, pp. 281-296.
- J. P. Paixão, I. M. Branco, E. Captivo, M. V. Pato, R. Eusébio and L. Amado (1986), "Bus and crew scheduling on a microcomputer", in *OR Models on Microcomputers*, J. D. Coelho and L. V. Tavares (eds.), North-Holland, pp.79-95.
- J. P. Paixão and M. V. Pato (1988), "Primal and dual greedy heuristics for the generalized set covering problem", *Investigação Operacional*, vol. 8, pp. 3-11.
- J. P. Paixão and M. V. Pato (1989), "A structural lagrangean relaxation for the two-duty period bus driver scheduling problem", *European Journal of Operational Research*, vol. 39, pp. 213-222.
- M. V. Pato (1989), *Algoritmos para Problemas de Cobertura Generalizados*, Ph. D. thesis, Faculdade de Ciências da Universidade de Lisboa.
- F. Shepardson and R. E. Marsten (1980), "A Lagrangean relaxation algorithm for the two duty period scheduling problem", *Management Science*, vol. 26, pp. 274-281.
- C. S. Wolfe (1984), "Cutting plane and branch and bound for solving a class of scheduling problems", *IIE Transactions*, vol. 16, pp. 50-58.
- L. Yihua (1985), "The application of the microcomputer in bus and crew scheduling in Shanghai", in *Computer Scheduling of Public Transport - 2*, J.-M. Rousseau (ed.), North-Holland, pp.179-198.



DIRECT METHODS FOR CONVEX QUADRATIC PROGRAMS SUBJECT TO BOX CONSTRAINTS

J.J. Júdice

Departamento de Matemática
Universidade de Coimbra
3000 Coimbra - Portugal

F.M. Pires

Departamento de Engenharia Civil
Universidade do Porto
4000 Porto - Portugal

Abstract: In this paper the authors investigate the use of direct methods for the solution of large-scale Convex Quadratic Programs subject to Box Constraints (BQP). An Active-Set Method, a Parametric Principal Pivoting Algorithm and a Bard-type Method are the algorithms considered in this study. Two implementations of these algorithms for the solution of large-scale Convex BQPs with general sparse and tridiagonal matrices are proposed. The implementations of the Active-Set Method can also be used to find stationary points (local minima) of Nonconvex BQPs.

The results of the experiences show that there is not a big gap between the efficiencies of the Active-Set and Bard-type methods but the former is consistently more efficient. The Parametric Principal Pivoting Method has the worst performance. Furthermore the implementations seem to be quite good in terms of storage, speed and accuracy.

Resumo: Neste artigo os autores investigam o uso de métodos directos para a solução de programas quadráticos de grandes dimensões com apenas restrições nos valores das variáveis (BQP). Um Método de Restrições Activas, um Algoritmo Pivotal Principal Paramétrico e um Método tipo Bard são os algoritmos testados neste estudo. São propostas duas implementações destes algoritmos para a resolução de BQPs Convexos de grandes dimensões com matrizes esparsas gerais e tridiagonais. As implementações do Método das Restrições Activas podem ser usadas para determinar pontos estacionários (mínimos locais) de BQPs não convexos.

Os resultados das experiências computacionais mostram que não é muito grande a diferença entre as eficiências dos métodos das Restrições Activas e tipo Bard, mas o primeiro algoritmo é sempre mais eficiente. O método Pivotal Principal Paramétrico é o menos eficiente. Além disso as implementações parecem ser muito boas em termos de armazenagem, rapidez de execução e precisão numérica.

Keywords: Quadratic Programming, Box Constraints, Linear Complementarity Problem, Large-Scale Problems, Sparse Matrices.

1 - Introduction

A Quadratic Program subject to Box Constraints (BQP) can be defined as

$$\begin{aligned} \text{Minimize} \quad & q^T z + \frac{1}{2} z^T M z \\ \text{subject to} \quad & a_i \leq z_i \leq b_i, \quad i = 1, 2, \dots, n \end{aligned} \tag{1}$$

where $q \in \mathbb{R}^n$, M is a symmetric matrix of order n and a_i, b_i are real numbers or $\pm \infty$ satisfying $a_i < b_i$. The BQP is Convex if M is a symmetric positive definite (SPD) or a singular symmetric positive semi-definite (SPSD) matrix. This problem has been studied during the past years and a large number of applications of the Convex BQP has been proposed. These include the solution of partial differential equations arising in Dirichlet problems with obstacles [13] or free-boundary value problems, such as the journal bearing [2, 23] and the elastic beam bending [39] problems. Optimal Control Problems [6] and Linear Least Squares [22, 25] are another two sources of Convex BQPs applications. A large number of applications has also been proposed in elastoplastic analysis of structures [29, 33], portfolio selection problems [32, 33] and spatial equilibrium models [34]. Convex BQPs are also useful in the solution of more general Convex Quadratic Programs [24]. Convex Nonlinear Programs subject to box constraints can be solved by SQP algorithms that rely on the solution of Convex BQPs [12]. In most of these applications the matrix M is large and sparse. This has increased the use of iterative methods to solve such problems [2, 6, 23, 24, 30, 40]. In this paper we study the applicability of direct methods to solve Convex BQPs.

Finding a global minimum for the Convex BQP(1) or a stationary point for the Nonconvex BQP(1) is equivalent solving the following problem

$$\begin{aligned} w &= q + Mz \\ a_i &\leq z_i \leq b_i \\ z_i = a_i &\Rightarrow w_i \geq 0 \\ z_i = b_i &\Rightarrow w_i \leq 0 \\ a_i < z_i < b_i &\Rightarrow w_i = 0 \end{aligned} \quad \left. \begin{array}{l} \\ \\ \\ \\ \end{array} \right\} \quad i = 1, 2, \dots, n \tag{2}$$

which represents its Kuhn-Tucker conditions [20]. The problem (2) can be seen as a generalization of the Linear Complementarity Problem(LCP)

$$w = q + Mz, \quad z \geq 0, \quad w \geq 0, \quad z^T w = 0 \tag{3}$$

Recently we have performed a computational study of the main direct methods for solving LCPs with SPD and singular SPSP matrices [18]. Murty's Bard-type method [28], Keller's method [19] and Graves' algorithm [14] have been considered in that study. We have proposed a static implementation for these methods that exploits the symmetry and sparsity of the matrix M . The study has shown that the efficiency of the first two algorithms is not very different, but usually Murty's method is the more efficient. Graves' method has the worst performance.

In this paper we present a similar computational study for the solution of large-scale Convex BQPs by direct methods. A Bard-type method [16], an Active-Set method [8] and a Parametric Principal Pivoting method [32] are the algorithms tested in this study. These three algorithms

can be seen as extensions of the algorithms referred to above. We introduce an implementation for these methods capable of solving large-scale Convex BQPs with general sparse SPD matrices. This procedure can be seen as a more advanced version of the implementation presented in [18]. We also propose another implementation for the solution of large-scale Convex BQPs with tridiagonal SPD matrices. The implementations of the Active-Set method can also be used to solve Convex BQPs with singular SPSP matrices and to find a stationary point (local minimum) of a Nonconvex BQP. The results of the computational experience show that:

- (i) The Bard-type and Active-Set methods usually perform in a similar way. However, the Active-Set method is consistently more efficient, even for the LCP(3). The Parametric Principal Pivoting method has in general the worst performance.
- (ii) The efficiency of the direct methods depends very much on the number nv of variables that are in both the initial and optimal active sets. The algorithms usually require a number of iterations slightly superior to $n-nv$, where n is the dimension of the Convex BQP, that is, the order of the matrix M . This number may be even worse when the global minimum is degenerate.
- (iii) Cycling may occur in the Bard-type algorithm but the phenomenon is quite rare.
- (iv) Both the implementations for general sparse and tridiagonal matrices seem to be quite good in terms of storage, speed and accuracy.

Active-Set methods require a feasible vector, that is a vector z satisfying the bounds, to start with. However, this is not necessary for the two remaining algorithms. We believe this may be exploited in the design of efficient procedures for solving Convex BQPs that do not depend so much on the starting vector. As is discussed in the last section of this paper, this may also have important implications on the design of Diagonalization and Sequential algorithms for Convex Nonlinear Programs.

There exist other procedures to solve Convex BQPs [37, 38]. However, it does not seem possible to design similar implementations for these methods that fully exploit the symmetry and sparsity of the matrix M . Hence they are not so recommended for solving large-scale Convex BQPs and are not discussed in this paper.

The organization of this paper is as follows. In Section 2 the direct methods are discussed. The implementations of these methods are described in Sections 3 and 4. Finally the computational experience with the direct methods is presented in the last section, together with some suggestions for future research.

2 - Direct Methods for Convex BQPs

Consider the BQP(1) and the problem (2). In order to describe the direct methods, we need to introduce the concept of Basic Solution for the problem (2). Let F be a subset of $\{1, \dots, n\}$ such that the principal submatrix M_{FF} of M is nonsingular. If

$$T = \{1, \dots, n\} - F = T_1 \cup T_2, \quad T_1 \cap T_2 = \emptyset$$

where \emptyset denotes the empty set, then $w = q + Mz$ is equivalent to

$$\begin{bmatrix} z_F \\ w_{T_1} \\ w_{T_2} \end{bmatrix} = \begin{bmatrix} \bar{q}_F \\ \bar{q}_{T_1} \\ \bar{q}_{T_2} \end{bmatrix} + \begin{bmatrix} \bar{M}_{FF} & \bar{M}_{FT_1} & \bar{M}_{FT_2} \\ \bar{M}_{T_1F} & \bar{M}_{T_1T_1} & \bar{M}_{T_1T_2} \\ \bar{M}_{T_2F} & \bar{M}_{T_2T_1} & \bar{M}_{T_2T_2} \end{bmatrix} \begin{bmatrix} w_F \\ z_{T_1} \\ z_{T_2} \end{bmatrix} \quad (4)$$

where

$$\bar{M}_{FF} = M_{FF}^{-1}, \bar{M}_{FT} = -M_{FF}^{-1} M_{FT}, \bar{M}_{TF} = M_{TF} M_{FF}^{-1}, \bar{q}_F = -M_{FF}^{-1} q_F \quad (5)$$

$$\bar{M}_{TT} = M_{TT} - M_{TF} M_{FF}^{-1} M_{FT}, \bar{q}_T = q_T - M_{TF} M_{FF}^{-1} q_F$$

A Basic Solution (\bar{z}, \bar{w}) for the problem (2) is a solution of the system (4) such that

$$F = \{i : w_i = 0\}, T_1 = \{i : z_i = b_i < +\infty\}, T_2 = \{i : z_i = a_i > -\infty\}$$

The variables $w_i, i \in F$ and $z_i, i \in T$ are called Nonbasic, while the remaining variables are said to be Basic. Hence the values of the nonbasic variables are

$$\bar{z}_{T_1} = b_{T_1}, \bar{z}_{T_2} = a_{T_2}, \bar{w}_F = 0 \quad (6)$$

Furthermore the formulas (4) and (5) imply that the values \bar{z}_F and \bar{w}_T of the basic variables satisfy

$$\begin{aligned} M_{FF} \bar{z}_F &= -(q_F + M_{FT_1} \bar{z}_{T_1} + M_{FT_2} \bar{z}_{T_2}) \\ \bar{w}_T &= q_T + M_{TT_1} \bar{z}_{T_1} + M_{TT_2} \bar{z}_{T_2} + M_{TF} \bar{z}_F \end{aligned} \quad (7)$$

where \bar{z}_{T_1} and \bar{z}_{T_2} are given by (6).

A Basic Solution (\bar{z}, \bar{w}) given by (6) and (7) is a solution of the problem (2) (\bar{z} is a global minimum of the Convex BQP or a stationary point of the Nonconvex BQP) if and only if

$$a_F \leq \bar{z}_F \leq b_F, \bar{w}_{T_1} \leq 0, \bar{w}_{T_2} \geq 0 \quad (8)$$

If this condition does not hold, then the basic solution is said to be Infeasible and one of the following cases may occur:

- (i) $\bar{z}_i < a_i$ ($i \in F$)
- (ii) $\bar{z}_i > b_i$ ($i \in F$)
- (iii) $\bar{z}_i = a_i$ and $\bar{w}_i < 0$
- (iv) $\bar{z}_i = b_i$ and $\bar{w}_i > 0$

which are called Infeasibilities of the problem (2).

The direct methods to be described in this paper are procedures that use in each iteration infeasible solutions of the problem (2) until finding a basic solution satisfying the condition (8). In each iteration of the algorithms an infeasibility is picked and the sets F, T_1 and T_2 are modified according to certain rules. Next, we describe the three direct methods.

(i) *Bard-Type Method*

As in Murty's Bard-type method [28], each iteration of this method consists of removing the first infeasibility associated with the current basic solution without caring about the feasibility of the remaining variables. The steps of the algorithm are as follows:

Step 0 - Let $F = \{i : a_i = -\infty \text{ and } b_i = +\infty\}$, $T_2 = \{i : a_i > -\infty\}$ and $T_1 = \{i : a_i = -\infty \text{ and } b_i < +\infty\}$

Step 1 - Compute \bar{z}_F and \bar{w}_T by (7) and let

$$I_1 = \{i \in F : \bar{z}_i > b_i\}, I_2 = \{i \in F : \bar{z}_i < a_i\} \tag{9}$$

$$I_3 = \{i \in T_1 : \bar{w}_i > 0\}, I_4 = \{i \in T_2 : \bar{w}_i < 0\}$$

and

$$s = \min \left\{ i \in \bigcup_{k=1}^4 I_k \right\} \tag{10}$$

If s does not exist, then $\bar{z} = (\bar{z}_F, b_{T_1}, a_{T_2})$ is a solution of the problem (2). Otherwise s belongs to exactly one of the four sets I_k . Go to step (k+1) provided $s \in I_k$.

Step 2 - $s \in I_1$ - Set $F = F - \{s\}$ and $T_1 = T_1 \cup \{s\}$ and go to Step 1.

Step 3 - $s \in I_2$ - Set $F = F - \{s\}$ and $T_2 = T_2 \cup \{s\}$ and go to Step 1.

Step 4 - $s \in I_3$ - If $a_s = -\infty$, set $F = F \cup \{s\}$ and $T_1 = T_1 - \{s\}$ and go to Step 1. Otherwise ($a_s > -\infty$) compute \bar{m}_{ss} and let

$$\theta = \min \left\{ \frac{\bar{w}_s}{\bar{m}_{ss}}, b_s - a_s \right\}$$

If $\theta = b_s - a_s$, set $T_1 = T_1 - \{s\}$ and $T_2 = T_2 \cup \{s\}$ and go to Step 1.

Otherwise set $T_1 = T_1 - \{s\}$ and $F = F \cup \{s\}$ and go to Step 1.

Step 5 - $s \in I_4$ - If $b_s = +\infty$, set $F = F \cup \{s\}$ and $T_2 = T_2 - \{s\}$ and go to Step 1.

Otherwise ($b_s < +\infty$) compute \bar{m}_{ss} and let

$$\theta = \min \left\{ -\frac{\bar{w}_s}{\bar{m}_{ss}}, b_s - a_s \right\}$$

If $\theta = b_s - a_s$, set $T_2 = T_2 - \{s\}$ and $T_1 = T_1 \cup \{s\}$ and go to Step 1.

Otherwise set $T_2 = T_2 - \{s\}$ and $F = F \cup \{s\}$ and go to Step 1.

It is easy to see that if $M \in \text{SPD}$, then all the four steps 2 to 5 can be performed and the infeasibility s is removed in all the cases. Furthermore, if M is a nonsingular M-matrix [35], then the algorithm converges to the unique solution of the problem (2) in a finite number of

iterations [16]. Hence the Bard-type method converges to the unique global minimum of the BQP if M is a SPD matrix with nonpositive off-diagonal elements. If M is a general SPD matrix, then it is not possible to prove convergence and we provide in appendix a Convex BQP for which cycling has occurred in the algorithm. However, cycling has occurred quite rarely in all the tests performed so far.

For a complete description of the algorithm we have to explain how the quantity \bar{m}_{ss} is computed. Since \bar{m}_{ss} is only necessary when $s \in T$, then, by using (5), we obtain the following procedure to compute this quantity

$$\begin{array}{l} \text{Solve } M_{FF} x = -M_{Fs} \rightarrow \bar{M}_{Fs} \\ \text{Compute } \bar{m}_{ss} = m_{ss} + M_{Fs}^T \bar{M}_{Fs} \end{array} \quad (11)$$

It follows from the description of the algorithm that it is necessary in each iteration to solve a system to compute the vectors \bar{z}_F and \bar{w}_T in Step 1. It may seem that if $s \in I_3$ or $s \in I_4$, then two systems have to be solved. Fortunately this is not true, since in these latter cases the values \bar{z}_F and \bar{w}_T in the next iteration can be updated by the following procedure:

$$\begin{array}{l} \bar{z}_s = \bar{z}_s - \varepsilon \theta, \quad \bar{z}_F = \bar{z}_F - \varepsilon \theta \bar{M}_{Fs} \\ \bar{w}_T = \bar{w}_T - \varepsilon \theta (M_{TF} \bar{M}_{Fs} + M_{Ts}) \end{array} \quad (12)$$

where θ is the quantity given in Step 4 or 5 and

$$\varepsilon = \begin{cases} 1 & \text{if } s \in T_1 \\ -1 & \text{if } s \in T_2 \end{cases}$$

These formulas can be obtained by linear algebra manipulations.

As in [18] we have replaced criterion (10) by one that chooses s as the row at which the largest infeasibility is achieved. Hence this new criterion consists of finding s as the first row at which

$$\max \left\{ \max \{ \bar{z}_i - b_i : i \in I_1 \}, \max \{ a_i - \bar{z}_i : i \in I_2 \}, \max \{ |\bar{w}_i| : i \in \bigcup_{k=3}^4 I_k \} \right\} \quad (14)$$

is achieved. We have not been able to prove convergence for this new version even for nonsingular M -matrices. Computational experience has shown that cycling is even rarer in this case.

(ii) Active-Set Method

This algorithm seeks a global minimum of the Convex BQP by only using solutions \bar{z} satisfying the bounds. Hence $I_1 = I_2 = \emptyset$ in any iteration of this algorithm, where I_k are the sets

given by (9). Therefore $(\bar{z}_F, a_{T_2}, b_{T_1})$ is a global minimum provided $I_3 = I_4 = \emptyset$. Otherwise let s be given by

$$s = \min \{ i : i \in I_3 \cup I_4 \} \tag{15}$$

or the first index at which the quantity

$$|\bar{w}_s| = \max \{ |\bar{w}_i| : i \in I_3 \cup I_4 \} \tag{16}$$

is attained. Then the value of the variable z_s is decreased (if $s \in T_1$) or increased (if $s \in T_2$) from its corresponding bound. The increase or decrease of the entering variable z_s may be unblocked or blocked by either the variable w_s or by a basic z_r ($r \in F$) whose value attains one of the bounds. If z_s is not blocked, then either attains the other bound or, in the case this bound is infinite, the quadratic function is unbounded below on its constraint set. If z_s is blocked by w_s , then z_s becomes basic by changing with w_s . If z_s is blocked by a variable z_r , $r \in F$, that attains a bound, then this latter variable becomes nonbasic and takes a value equal to that bound. A cycle is initiated in which z_s stays nonbasic with a value \bar{z}_s between the bounds and continues to be increased or decreased until one of the previous cases occurs. During this cycle only a finite number of basic variables z_i ($i \in F$) becomes nonbasic.

The steps of the Active-Set method are presented below.

Step 0 - Let $z = (\bar{z}_F, b_{T_1}, a_{T_2})$ be an initial solution such that $a_F \leq \bar{z}_F \leq b_F$ and compute \bar{w}_T by (7).

Step 1 - If $I_3 = I_4 = \emptyset$, then $(\bar{z}_F, b_{T_1}, a_{T_2})$ is a global minimum of the Convex BQP and stop. Otherwise let s be the index given by (15) or (16).

Step 2 - Compute \bar{M}_{FS} and \bar{m}_{ss} by (11) and let ϵ be given by (13). Compute θ_1 , θ_2 and θ_3 by

$$\theta_1 = \begin{cases} \frac{\bar{w}_s}{\epsilon \bar{m}_{ss}} & \text{if } \bar{m}_{ss} > 0 \\ +\infty & \text{if } \bar{m}_{ss} \leq 0 \end{cases}$$

$$\theta_2 = \begin{cases} \frac{\bar{z}_r - a_r}{\epsilon \bar{m}_{rs}} = \min \left\{ \frac{\bar{z}_i - a_i}{\epsilon \bar{m}_{is}} : i \in F \text{ and } \epsilon \bar{m}_{is} > 0 \right\} \\ +\infty & \text{if } F = \emptyset \text{ or } \epsilon \bar{m}_{is} \leq 0 \text{ for all } i \in F \end{cases}$$

$$\theta_3 = \begin{cases} \frac{b_r - \bar{z}_r}{-\epsilon \bar{m}_{rs}} = \min \left\{ \frac{b_i - \bar{z}_i}{-\epsilon \bar{m}_{is}} : i \in F \text{ and } \epsilon \bar{m}_{is} < 0 \right\} \\ +\infty & \text{if } F = \emptyset \text{ or } \epsilon \bar{m}_{is} \geq 0 \text{ for all } i \in F \end{cases}$$

and let $\theta = \min \{ b_s - a_s, \theta_1, \theta_2, \theta_3 \}$. Go to Step (k+2) provided $s \in T_k$, $k=1, 2$.

Step 3 - $s \in T_1$ -

- (i) If $\theta = +\infty$, the quadratic function is unbounded below on its constraint set and stop.
- (ii) If $\theta = b_s - a_s$, set $T_1 = T_1 - \{s\}$ and $T_2 = T_2 \cup \{s\}$. Update \bar{z} and \bar{w} by (12) and go to Step 1.
- (iii) If $\theta = \theta_1$, set $T_1 = T_1 - \{s\}$ and $F = F \cup \{s\}$. Update \bar{z} and \bar{w} by (12) and go to Step 1.
- (iv) If $\theta = \theta_2$, set $T_2 = T_2 \cup \{r\}$ and $F = F - \{r\}$. Set $\bar{z}_r = a_r$ and update \bar{z}_F, \bar{z}_s and \bar{w}_s by (12). Go to Step 2.
- (v) If $\theta = \theta_3$, set $T_1 = T_1 \cup \{r\}$ and $F = F - \{r\}$. Set $\bar{z}_r = b_r$ and update \bar{z}_F, \bar{z}_s and \bar{w}_s by (12). Go to Step 2.

Step 4 - $s \in T_2$ -

- (i) If $\theta = +\infty$, the quadratic function is unbounded below on its constraint set and stop.
- (ii) If $\theta = b_s - a_s$, set $T_2 = T_2 - \{s\}$ and $T_1 = T_1 \cup \{s\}$. Update \bar{z} and \bar{w} by (12) and go to Step 1.
- (iii) If $\theta = \theta_1$, set $T_2 = T_2 - \{s\}$ and $F = F \cup \{s\}$. Update \bar{z} and \bar{w} by (12) and go to Step 1.
- (iv) If $\theta = \theta_2$, set $T_2 = T_2 \cup \{r\}$ and $F = F - \{r\}$. Set $\bar{z}_r = a_r$ and update \bar{z}_F, \bar{z}_s and \bar{w}_s by (12). Go to Step 2.
- (v) If $\theta = \theta_3$, set $T_1 = T_1 \cup \{r\}$ and $F = F - \{r\}$. Set $\bar{z}_r = b_r$ and update \bar{z}_F, \bar{z}_s and \bar{w}_s by (12). Go to Step 2.

To complete the description of the Active-Set method, it is only necessary to explain the process of finding an initial solution in Step 0. Consider the sets

$$G = \{ i : a_i = -\infty \text{ and } b_i = +\infty \} \quad (17)$$

$$T_2 = \{ i : a_i > -\infty \}, \quad T_1 = \{ i : a_i = -\infty \text{ and } b_i < +\infty \}$$

Then two cases may occur:

- (i) If $M_{GG} \in \text{SPD}$, then set $F = G$ and let \bar{z}_F be the solution of the following linear system

$$M_{FF} x = - (q_F + M_{FT_1} b_{T_1} + M_{FT_2} a_{T_2}) \quad (18)$$

Then $(\bar{z}_F, b_{T_1}, a_{T_2})$ is the initial solution.

- (ii) If M_{GG} is a singular SPSD matrix then [18] either the quadratic function is unbounded on its constraint set or there exists a set $\emptyset \subsetneq F \subset G$ such that $M_{FF} \in \text{SPD}$. In this latter case all the variables $z_i, i \in G - F$, can be set to zero until the end of the procedure, that is, we set $a_i = b_i = 0$ for $i \in G - F$ and $T_2 = T_2 \cup (G - F)$. The initial solution is $(\bar{z}_F, a_{T_2}, b_{T_1})$, where \bar{z}_F is the solution of the linear system (18).

The description of the algorithm shows that it is required in each iteration to solve a linear system with matrix M_{FF} and perform some operations with matrices and vectors, scalar products of vectors and sums of vectors. Moreover M_{FF} is initially a SPD matrix (we assume that the determinant $\det(M_{GG})$ of M_{GG} is equal to one) and only increases when $\bar{m}_{ss} > 0$. In this last case we have [3]

$$\det \begin{pmatrix} M_{FF} & M_{Fs} \\ M_{sF} & m_{ss} \end{pmatrix} = \det(M_{FF}) \bar{m}_{ss} > 0$$

Hence all the submatrices M_{FF} used by the algorithm are SPD and the Active-Set method can be used to solve any Convex BQP.

It is possible to show by using linear algebra manipulations that the value of the quadratic function reduces in each iteration. Hence the Active-Set method always terminates in a finite number of iterations and either finds a global minimum of a Convex BQP or shows that the quadratic function is unbounded below on its constraint set. This latter case cannot occur if $M \in \text{SPD}$.

It is important to mention that this algorithm is exactly the Fletcher-Jackson Active-Set method [8] presented in a different form. We have chosen this new presentation in order to better see the similarities and differences among the three direct methods discussed in this paper. The Active-Set method can be seen as a special case of the Bard-type method in which the sets I_1 and I_2 are forced to be empty in each iteration.

The Active-Set Method can be used to find a local minimum of a Concave BQP, that is, a BQP in which $-M \in \text{SPSD}$. If G is the set given by (17), then three cases may occur:

- (i) If $G = \emptyset$, then the Active-Set method is used without any modification.
- (ii) If $G \neq \emptyset$, and $q_G = 0$, $M_{GG} = 0$, then set $a_i = b_i = 0$ for $i \in G$ and $T_2 = T_2 \cup G$ in the Initial Step, and use the Active-Set method as before.
- (iii) If the two previous cases do not occur, the quadratic function is unbounded from below and the algorithm stops.

It is worth mentioning that for Concave BQPs the set F is always empty in any iteration of the algorithm. This means that no systems have to be solved in this case. Furthermore the Bard-type method coincides with the Active-Set method for this kind of BQP.

The Active-Set method can also be used to find a stationary point of a Nonconvex BQP (that is, a BQP such that neither $M \in \text{SPSD}$ nor $-M \in \text{SPSD}$) provided the set G given by (17) is empty or $M_{GG} \in \text{SPD}$. If none of these conditions holds, then it is necessary to make a change of variables from $z_i, i \in G$ into z'_i by the transformation

$$z_i = z'_i - z_{n+1}, \quad i \in G$$

and write $a_{n+1} = a_i = 0, i \in G$. Then $G = \emptyset$ for the resulting BQP and the Active-Set method can be used to find a stationary point for this program. This stationary point is a local minimum provided $w_{T_1} < 0$ and $w_{T_2} > 0$. If $w_{T_1} \leq 0$ and $w_{T_2} \geq 0$ and there exists at least a component $w_i = 0 (i \in T)$ then the stationary point is usually a local minimum but it can also be a saddle-point [12].

(iii) *Parametric Principal Pivoting Method*

This algorithm is essentially due to Pang [32] and can be seen as an extension of Graves' method [14] for the solution of Convex BQPs with $M \in \text{SPD}$. As in this method, an additional vector $p \in \mathbb{R}^n$ and a parameter λ are introduced in $w = q + Mz$ to get a system of the form $w = q + \lambda p + Mz$. The initial solution

$$(\bar{z}_F, b_{T_1}, a_{T_2}) \quad (19)$$

is obtained as in the Bard-type method and p is any vector satisfying

$$\begin{cases} p_i > 0 & \text{if } i \in T_2 \\ p_i < 0 & \text{if } i \in T_1 \\ p_i = 0 & \text{if } i \in F \end{cases}$$

where F , T_1 and T_2 are the sets defined in Step 0 of the Bard-type method. If (19) is the initial solution, then there exists a value $\bar{\lambda} > 0$ such that

$$\bar{z}_F = \bar{z}_F + \bar{\lambda} \bar{p}_F, \quad \bar{z}_{T_1} = b_{T_1}, \quad \bar{z}_{T_2} = a_{T_2} \quad (20)$$

is the global minimum of the Convex BQP_{λ} given by

$$\begin{aligned} \text{Minimize } & (q + \lambda p)^T z + \frac{1}{2} z^T M z \\ \text{subject to } & a_i \leq z_i \leq b_i, \quad i = 1, 2, \dots, n \end{aligned} \quad (21)$$

The Parametric Principal Pivoting Method consists of finding a nonincreasing sequence of values $\bar{\lambda}$ of λ such that the solution (20) is a global minimum of $\text{BQP}_{\bar{\lambda}}$. A global minimum of the BQP_0 is a global minimum of the Convex $\text{BQP}(1)$. As in parametric linear programming, the vector p is transformed in each iteration. By using similar formulas to those presented in (5), we can easily see that in any iteration the corresponding vector \bar{p} is obtained by the following procedure:

Solve $M_{FF} x = -p_F \rightarrow \bar{p}_F$ Compute $\bar{p}_T = p_T + M_{TF} \bar{p}_F$	(22)
---	------

The steps of the Principal Pivoting Method are as follows:

Step 0 - Let $F = \{ i : a_i = -\infty \text{ and } b_i = +\infty \}$, $T_2 = \{ i : a_i > -\infty \}$ and $T_1 = \{ i : a_i = -\infty \text{ and } b_i < +\infty \}$.

General Step - Compute \bar{z}_F, \bar{w}_T by (7) and \bar{p} by (22). Let I_k , $k = 1, \dots, 4$ be the sets of infeasibilities given by (9). If

$$\bigcup_{k=1}^4 I_k = \emptyset$$

then $(\bar{z}_F, a_{T_2}, b_{T_1})$ is the global minimum of the BQP and stop.

Otherwise compute

$$\lambda_1 = \max \left\{ \frac{b_i - \bar{z}_i}{p_i} : i \in I_1 \right\}$$

$$\lambda_2 = \max \left\{ \frac{a_i - \bar{z}_i}{p_i} : i \in I_2 \right\}$$

$$\lambda_3 = \max \left\{ -\frac{\bar{w}_i}{p_i} : i \in I_3 \cup I_4 \right\}$$

Let $\bar{\lambda} = \max \{ \lambda_1, \lambda_2, \lambda_3 \}$ and s be the first index at which the value $\bar{\lambda}$ is attained. Then:

- (i) if $s \in I_1$, set $F = F - \{s\}$ and $T_1 = T_1 \cup \{s\}$
- (ii) if $s \in I_2$, set $F = F - \{s\}$ and $T_2 = T_2 \cup \{s\}$
- (iii) if $s \in I_3$, set $F = F \cup \{s\}$ and $T_2 = T_2 - \{s\}$
- (iv) if $s \in I_4$, set $F = F \cup \{s\}$ and $T_1 = T_1 - \{s\}$

Repeat General Step.

It follows from the description of the algorithm that two systems with matrix M_{FF} have to be solved in each iteration of the algorithm. This is an obvious disadvantage of the algorithm. The method terminates in a finite number of iterations with the unique global minimum of the Convex BQP provided a tie-break rule similar to [27, Chapter 11] is incorporated to avoid cycling in degenerate cases. However, the choice of the first index in case of ties has proven sufficient to guarantee convergence in all the tests performed so far.

3 - Implementation of the Algorithms for a General Sparse Matrix

In this section we describe an implementation for the three algorithms discussed in the previous section when M is a sparse SPD or a singular SPSD matrix without special structure. The procedure borrows some ideas from the solution of large-scale linear systems with such matrices [7].

(i) *ANALYSE Phase and Data Structures for the matrices M and LDL^T decomposition of M_{FF}*

As discussed in the previous section, it is necessary in each iteration of the algorithms to solve linear systems with matrix M_{FF} . The BQP is invariant under symmetric permutations. In addition either the quadratic function is unbounded below on its constraint set (this can only occur when $M \notin \text{SPD}$) or it is always possible to remove the rows and columns s such that $m_{ss} = 0$ [18]. So we can assume that all the diagonal elements of M are positive. These

properties enable the use of an ordering algorithm, such as the minimum-degree algorithm to find a symmetric permutation for the rows and columns of the matrix M that is suitable in sparse terms [11]. So, instead of solving the BQP, we solve the permuted BQP that arises from this ordering algorithm.

After the minimum-degree algorithm is applied, the data structure for the permuted matrix M is created. This matrix is stored columnwise as an unsymmetric matrix, in order to facilitate the matrix-vector products given by (7), (11), (12) and (22) and also the updating procedures for the LDL^T decompositions of the submatrices M_{FF} to be described later. Hence the matrix M is stored by using the arrays DM , IM , PM and VM stated below.

DM - array of dimension n containing all the diagonal elements of M .

PM - array of dimension $(n+1)$ containing pointers for the first off-diagonal entry in each column of the matrix M . $PM(n+1)$ points to the position after the last nonzero entry of M . Note that $PM(j+1) - PM(j)$ is the number of nonzero elements in column j .

IM - array of dimension equal to the number of nonzero off-diagonal elements of M containing the row indices corresponding to each nonzero off-diagonal entry of M .

VM - array of dimension equal to the number of nonzero off-diagonal elements of M containing the values of these elements.

To clarify the understanding of this data structure, let the permuted matrix M be given by

$$M = \begin{bmatrix} 1. & 0. & 0. & 1. & 0. \\ 0. & 3. & 2. & 1. & 0. \\ 0. & 2. & 2. & 0. & 0. \\ 1. & 1. & 0. & 5. & 3. \\ 0. & 0. & 0. & 3. & 4. \end{bmatrix} \quad (23)$$

Then

$$\begin{aligned} DM &= \left| \begin{array}{c} 1. \\ 3. \\ 2. \\ 5. \\ 4. \end{array} \right| & IM &= \left| \begin{array}{c} 4 \\ 3 \\ 4 \\ 2 \\ 1 \\ 2 \\ 5 \\ 4 \end{array} \right| \\ PM &= \left| \begin{array}{c} 1 \\ 2 \\ 4 \\ 5 \\ 8 \\ 9 \end{array} \right| & VM &= \left| \begin{array}{c} 1. \\ 2. \\ 1. \\ 2. \\ 1. \\ 1. \\ 3. \\ 3. \end{array} \right| \end{aligned}$$

Another data structure is necessary to store all the matrices L and D corresponding to the LDL^T decompositions of the matrices M_{FF} used by the algorithms. To create this data structure, the symbolic phase described in [11, Chapter 5] is applied. Let $nonz$ and $nfill$ be respectively the number of off-diagonal elements of M strictly below the diagonal and the fill-in positions that occur in the symbolic phase. After this procedure is applied it is known the maximum number of off-diagonal elements that each column of L can take. This allows us to construct an array of pointers PL of dimension n such that $PL(i)$ points to the first off-diagonal entry in the column i of the matrix L of the LDL^T decomposition of M . In the case of the matrix (23), we have $nfill = 1$ and $(4, 3)$ is the unique fill-in position, whence PL has the following form

$$PL = \left| \begin{array}{c} 1 \\ 2 \\ 4 \\ 5 \\ 6 \end{array} \right| \quad (24)$$

This array is never modified during the whole procedure. Besides this unmodified array, the

data structure for the matrices L and D of the LDL^T decompositions of the successive submatrices M_{FF} contains four other arrays that are modified in each iteration. These are presented below.

DD - array of dimension n containing all the diagonal elements of the matrix D.

IL - array of dimension (nonz + nfill) containing all the row indices corresponding to the nonzero off-diagonal elements of the matrix L.

VL - array of dimension (nonz + nfill) containing the values of the nonzero off-diagonal elements in each column of the matrix L.

NZL - array of dimension n containing the number of off-diagonal elements of the matrix L.

In each iteration the elements of column $i \in F$ of the matrix L of the LDL^T decomposition of M_{FF} are located in IL and VL between the components $PL(i)$ and $PL(i) + NZ(i) - 1$. The next components from $PL(i) + NZ(i)$ to $PL(i+1) - 1$ are ignored. Also ignored are the entries $IL(j)$ and $VL(j)$ such that $j \in \{PL(t), PL(t+1) - 1\}$ and $t \notin F$.

The diagonal elements of the matrix D of the LDL^T decomposition of M_{FF} are located in the positions $i \in F$ of DD. The components j of DD such that $j \notin F$ are ignored. Furthermore the components j of NZL such $j \in F$ are also ignored.

(ii) *The Initial Step*

The ANALYSE Phase is completed after the vector PL is constructed. Then storage space is provided for the remaining arrays of the data structure of the matrices L and D of the LDL^T decompositions of the matrices M_{FF} . The FACTOR-SOLVE Phase starts by finding the LDL^T decomposition of the matrix $M_{FF} \in SPD$, where $F = G$ and G is the set given by (17). The elements of the matrices L and D are placed in the respective positions according to the scheme explained above. To illustrate this step consider the Convex BQP with matrix $M \in SPD$ given by (23) and suppose that $G = \{2, 4\}$. As stated before, $nfill = 1$ and (4, 3) is the fill-in position for this matrix. Furthermore

$$M_{FF} = \begin{bmatrix} 3. & 1. \\ 1. & 5. \end{bmatrix} = \begin{bmatrix} 1. & \\ 1/3 & 1. \end{bmatrix} \begin{bmatrix} 3. & \\ & 14/3 \end{bmatrix} \begin{bmatrix} 1 & 1/3 \\ & 1. \end{bmatrix}$$

If - represents an ignored component, then

$$\begin{aligned} DD &= \begin{bmatrix} - & 3. & - & - & 14/3 & - \\ & & & & & \end{bmatrix} & \quad IM &= \begin{bmatrix} - & 4 & - & - & - \\ & & & & \end{bmatrix} \\ VL &= \begin{bmatrix} - & 1/3 & - & - & - \\ & & & & \end{bmatrix} & \quad NZL &= \begin{bmatrix} - & 1 & - & 0 & - \\ & & & & \end{bmatrix} \end{aligned} \tag{25}$$

which together with the array (24) constitute the data structure for the matrices L and D of the LDL^T decomposition of M_{FF} .

As stated in the previous section, the Initial Step is not performed if it is empty the set G given by (17). If $G \neq \emptyset$ but $M_{GG} \notin SPD$, this is noticed during the procedure of finding the LDL^T decomposition of the matrix M_{GG} by the occurrence of a zero diagonal element d_{ii} (d_{ii} smaller than a tolerance for zero) in the matrix D. In this last case \bar{w}_i is computed by using the formulas (7) and the LDL^T decomposition of the submatrix $M_{FF} \in SPD$ of M_{GG} already

obtained. Two cases may occur:

- (i) if $\bar{w}_i \neq 0$, the quadratic function is unbounded below on its constraint set and the algorithm stops.
- (ii) if $\bar{w}_i = 0$, then the row and column i are deleted from further considerations. As stated before we set $a_i = b_i = 0$ and $T_2 = T_2 \cup \{i\}$ and we proceed with the next index of G .

(iii) Data Structure for the sets F , T_1 and T_2

In each iteration of the direct methods the sets F , T_1 and T_2 are modified according to certain rules. Hence it is necessary to construct a data structure for the representation of these sets. This data structure should be designed in such a way that not only the sets are easily updated but also the matrix-vector products given in (7), (11) (12) and (22) are performed quickly. We propose a data structure consisting of two integer arrays IV and IVS of dimension n . If $|F|$, $|T_1|$ and $|T_2|$ are the numbers of elements of the sets F , T_1 and T_2 respectively, then

- (i) the first $|F|$ components of IV contain the indices of F
- (ii) the last $|T_1|$ components of IV contain the indices of T_1
- (iii) the remaining components of IV contain the indices of T_2

Each component i of the array IVS is positive or negative depending on $i \in F$ or $i \in T_1 \cup T_2$. Furthermore the absolute value of IVS (i) represents the position in IV where the index i can be found. For instance for the BQP with matrix M given by (23) the data structure corresponding to a solution such that $T_1 = \{3\}$, $F = \{2, 4\}$ and $T_2 = \{1, 5\}$ is as follows

$$IV = \left[\begin{array}{cccccc} 2 & 4 & 1 & 5 & 3 & \end{array} \right] \quad IVS = \left[\begin{array}{cccccc} -3 & 1 & -5 & 2 & -4 & \end{array} \right]$$

It is not difficult to see that both the updating of the sets F , T_1 and T_2 and the evaluation of the matrix-vector products required by the algorithms can be performed quickly by using this data structure.

(iv) Procedure for updating the LDL^T decomposition of M_{FF}

Whenever the set F is modified the LDL^T decomposition of the resulting matrix M_{FF} must be updated from the LDL^T decomposition of the old matrix M_{FF} . The updating must be performed by following the ordering achieved in ANALYSE Phase. To explain the updating procedure, consider first the case where the set F is increased by an index s . Let the LDL^T decomposition of the old matrix M_{FF} be given by

$$M_{FF} = \begin{bmatrix} A & B^T \\ B & C \end{bmatrix} = \begin{bmatrix} L_1 & \\ & L_2 \end{bmatrix} \begin{bmatrix} D_1 & \\ & D_2 \end{bmatrix} \begin{bmatrix} L_1^T & E^T \\ & L_2^T \end{bmatrix} \tag{26}$$

where L_i and D_i , $i = 1, 2$ are unit lower triangular and diagonal matrices respectively. If $\bar{F} = F \cup \{s\}$ is the resulting set, then we can write

$$M_{\bar{F}\bar{F}} = \begin{bmatrix} A & d & B^T \\ d^T & \alpha_0 & e^T \\ B & e & C \end{bmatrix} \begin{matrix} \leftarrow \text{row } s \\ \uparrow \\ \text{column } s \end{matrix} \tag{27}$$

where α_0 is a real number and d, e are vectors of appropriate dimensions. Then [18]

$$M_{\overline{F}\overline{F}} = \begin{bmatrix} L_1 & & & \\ \overline{d}^T & 1 & & \\ E & \overline{e} & \overline{L}_2 & \end{bmatrix} \begin{bmatrix} D_1 & & & \\ & \overline{\alpha}_0 & & \\ & & D_2 & \\ & & & \end{bmatrix} \begin{bmatrix} L_1^T & \overline{d} & E^T \\ & 1 & \overline{e}^T \\ & & \overline{L}_2^T \end{bmatrix} \quad (28)$$

where the real number $\overline{\alpha}_0$ and the vectors \overline{d} and \overline{e} are given by the following procedure

Solve the system $L_1 x = d$

$$\text{Compute } \begin{cases} \overline{\alpha}_0 = \alpha_0 - x D_1^{-1} x \\ \overline{e} = \frac{1}{\overline{\alpha}_0} (e - E x) \\ \overline{d} = D_1^{-1} x \end{cases} \quad (29)$$

and the matrices \overline{L}_2 and \overline{D}_2 are the unit lower triangular and the diagonal matrices of the LDL^T decomposition of the matrix

$$L_2 D_2 L_2^T - \overline{\alpha}_0 \overline{e} \overline{e}^T$$

This decomposition can be obtained by a version for sparse matrices [21] of Bennett's algorithm [1]. This last procedure incorporates a symbolic phase that indicates the elements that have to be modified in each column. Care must be taken with elements that become zero due to cancellation and are considered to be nonzero by this symbolic phase. They must be stored with value zero. We also note that Bennett's algorithm is not used if the index s to be added to F satisfies $s \geq j$ for all $j \in F$ in the ordering achieved in ANALYSE Phase.

The procedure for updating the LDL^T decomposition when an index s is deleted from the set F is even simpler. Let \overline{F} be the old set and $F = \overline{F} - \{s\}$. If M_{FF} and $M_{\overline{F}\overline{F}}$ are the corresponding matrices given by (26) and (27) then it is only necessary to compute the LDL^T decomposition of the matrix

$$\overline{L}_2 \overline{D}_2 \overline{L}_2^T + \alpha_0 \overline{e} \overline{e}^T$$

which can be done by the sparse version of Bennett's algorithm stated above. Note that no work has to be done if s is the last row and column of the matrix $M_{\overline{F}\overline{F}}$.

By using some contrived examples, Fletcher and Powell [9] have shown that when s is added to F Bennett's algorithm may have some difficulties in presence of rounding errors. They have designed a more elaborate algorithm to overcome these problems. In practice, these difficulties are checked by the occurrence of a nonpositive element (smaller than a tolerance for zero) in the matrix D of the LDL^T decomposition of the matrix M_{FF} . Our experience has shown that this phenomenon is extremely rare. So instead of using Fletcher and Powell algorithm, we have decided to implement Bennett's method and use a reinversion whenever a nonpositive diagonal element of a matrix D occurs. This reinversion consists of finding directly the LDL^T decomposition of M_{FF} from stack instead of using the updating procedure.

Whenever an updating of the LDL^T decomposition is performed, the data structure for the matrices L and D has to be modified. However, it is important to stress that not all the components of this data structure are modified but only those corresponding to the columns and rows that are active in the symbolic phase of this updating. To illustrate how this data structure is modified, consider again the BQP whose matrix M is given by (23). Suppose that F = {2, 4} and let the data structure for the matrices L and D of the decomposition of M_{FF} be given by (25). Furthermore let 1 the index to be added to the set F. Then F becomes {1, 2, 4}. Since

$$M_{FF} = \begin{bmatrix} 1. & 0. & 1. \\ 0. & 3. & 1. \\ 1. & 1. & 5. \end{bmatrix} = \begin{bmatrix} 1. & & \\ 0. & 1. & \\ 1. & 1/3 & 1. \end{bmatrix} \begin{bmatrix} 1. & & \\ & 3. & \\ & & 14/3 \end{bmatrix} \begin{bmatrix} 1. & 0. & 1. \\ & 1. & 1/3 \\ & & 1. \end{bmatrix}$$

then the data structure for the resulting matrices is given by the PL unmodified array and

$$DD = \left| \begin{array}{cccc} 1. & 3. & - & 14/3 & - \end{array} \right| \quad IL = \left| \begin{array}{cccc} 4 & 4 & - & - & - \end{array} \right|$$

$$VL = \left| \begin{array}{cccc} 1. & 1/3 & - & - & - \end{array} \right| \quad NZL = \left| \begin{array}{cccc} 1 & 1 & - & 0 & - \end{array} \right|$$

To complete the description of the implementation, we should add that the linear systems with matrix M_{FF} are solved by using the data structure that stores the matrices L and D of the LDL^T decomposition of this matrix and the algorithms described in [11]. The matrix-vector products and scalar products stated in (7), (11), (12) and (22) are performed by using the data structure for the matrix M, the data structure for the index sets F, T₁ and T₂ and the vector \bar{z}_F or \bar{M}_{F_S} solution of the linear systems with matrix M_{FF}.

We have described an implementation for the three direct methods that exploits the symmetry and sparsity of the matrix of M of the Convex BQP. The procedure is static, since the dimensions are not modified and garbage collections [31] are not necessary during the whole procedure. The total primary storage for this static implementation is given by

$$\text{stor} = 2 \times n + 3 \times \text{nonz} + \text{nfill} \quad (30)$$

where nonz and nfill are defined as before. Therefore this implementation is recommended whenever there is sufficient storage to accommodate stor real numbers and a similar number of integer numbers that are necessary for the data structures of the index sets and matrices L and M.

(v) Implementation for a Nonconvex BQP

Consider now the case of a Nonconvex BQP. Since M_{FF} ∈ SPD in each iteration of the Active-Set method then only the principal submatrix of M with positive diagonal elements should be considered in the ANALYSE Phase. This means that the value of stor is reduced for Nonconvex BQPs. Furthermore if it is empty the set G defined by (17) or M_{GG} ∈ SPD, then the procedure described in this section with the modification stated above is a valid implementation for the Active-Set Method. If G ≠ ∅ and M_{GG} ∉ SPD, then as stated before it is required to make the change of variables

$$z_G = z_G - z_{n+1} e_G$$

where e_G is a vector of ones of order $|G|$. If $H = \{1, \dots, n\} - G$, then the BQP(1) is replaced by the equivalent Nonconvex BQP

$$\text{Min} \begin{bmatrix} q_G \\ q_H \\ q_G^T e_G \end{bmatrix}^T \begin{bmatrix} z_G \\ z_H \\ z_{n+1} \end{bmatrix} + \frac{1}{2} \begin{bmatrix} z_G \\ z_H \\ z_{n+1} \end{bmatrix}^T \begin{bmatrix} M_{GG} & M_{HG}^T & M_{GG} e_G \\ M_{HG} & M_{HH} & M_{HG} e_G \\ e_G^T M_{GG} & e_G^T M_{HG}^T & e_G^T M_{GG} e_G \end{bmatrix} \begin{bmatrix} z_G \\ z_H \\ z_{n+1} \end{bmatrix}$$

subject to $a_i \leq z_i \leq b_i, i = 1, 2, \dots, n+1$

where $a_i = 0, b_i = +\infty$ for $i \in G$ and $a_{n+1} = 0, b_{n+1} = +\infty$. Therefore the matrix of this BQP has the form

$$\begin{bmatrix} M & f \\ f^T & \beta_0 \end{bmatrix}$$

where β_0 is a real number, f is a dense vector of order n and M is the matrix of the BQP(1).

As before, in the implementation for such a Nonconvex BQP only the principal submatrix of M with positive diagonal elements is considered in the ANALYSE Phase. The data structures are defined as before and stor is increased by n , since the vector f also has to be stored. The matrix of the systems required by the Active-Set method either have the form M_{FF} with $F \subseteq \{1, \dots, n\}$ or

$$\begin{bmatrix} M_{FF} & f_F \\ f_F^T & \beta_0 \end{bmatrix}$$

depending on z_{n+1} being a nonbasic or a basic variable respectively. In the latter case the solution of a system of the form

$$\begin{bmatrix} M_{FF} & f_F \\ f_F^T & \beta_0 \end{bmatrix} \begin{bmatrix} x \\ y_0 \end{bmatrix} = \begin{bmatrix} \gamma \\ \gamma_0 \end{bmatrix}$$

is performed by the following procedure:

Solve	$M_{FF} x = \gamma$
Solve	$M_{FF} \beta = f_F$
Compute	$y_0 = \frac{\gamma_0 - f_F^T x}{\beta_0 - f_F^T \beta}$
Compute	$x = x - y_0 \beta$

These are the main features to be added to the implementation for the solution of Nonconvex BQPs when $G \neq \emptyset$ and $M_{GG} \in \text{SPD}$.

The implementation for a Concave BQP is even simpler. In fact, since no systems have to be solved for this kind of problem, then no ordering algorithm is required in ANALYSE Phase and the data structure for the matrices L and D is not necessary. This reduces the value of stor to $n + 2 \times \text{nonz}$ and allows the solution of quite large Concave BQPs.

4 - Implementation of the Algorithms for Tridiagonal Matrices

Tridiagonal matrices constitute an important example of structured matrices that occur quite frequently in applications. Due to this special structure, it is possible to design an implementation for the direct methods that is more efficient for these BQPs than the procedure explained in the previous section. Since M is a tridiagonal matrix, then it can be stored by only using an array D of dimension n that stores its diagonal elements and an array of dimension $(n-1)$ for its subdiagonal elements. Furthermore M_{FF} is a block diagonal matrix such that each diagonal block is itself a tridiagonal matrix (we can assume that matrices of order 1 and 2 are tridiagonal). Then we can write

$$M_{FF} = \text{diag} (M_{F_1 F_1}, \dots, M_{F_k F_k}) \quad (31)$$

where

$$F = \bigcup_{i=1}^k F_i$$

$$F_i \cap F_j = \emptyset \text{ for } i \neq j \quad (32)$$

$$\max \{ s : s \in F_i \} < \min \{ r : r \in F_{i+1} \} - 1 \text{ for each } i=1, \dots, k-1$$

For instance if $n = 10$ and $F = \{1, 2, 3, 6, 8, 9\}$ then

$$M_{FF} = \begin{bmatrix} m_{11} & m_{12} & & & & & & & & \\ & m_{21} & m_{22} & m_{23} & & & & & & \\ & & m_{32} & m_{33} & & & & & & \\ & & & & m_{66} & & & & & \\ & & & & & & m_{88} & m_{89} & & \\ & & & & & & m_{98} & m_{99} & & \end{bmatrix} \quad (33)$$

and we have the representation stated above with $F_1 = \{1, 2, 3\}$, $F_2 = \{6\}$, $F_3 = \{8, 9\}$.

As suggested in [5], each matrix M_{FF} can be recovered from M by using a linked list consisting of a scalar IP and an array of dimension IB of dimension $(n+1)$. The scalar IP and the array IB are defined as follows:

IP = first row and column of the first diagonal block of M_{FF} .

$IB(i)$ = first row and column of the next diagonal block (if $IB(i) = 0$ then the current block is the last).

$IB(i+1)$ = number of rows (and columns) of the diagonal block that starts at the row (and column) i .

For instance the linked list associated with the matrix (33) is given by

$$IP = 1, \quad IB = (6, 3, -, -, -, 8, 1, 0, 2, -, -)$$

where $-$ means an ignored location.

In each iteration of the direct methods the set F is modified in at most one element. A change in the set F implies a modification in the diagonal blocks of the matrix M_{FF} , that is, in the linked list. For instance, if the index 2 is taken from F , then M_{FF} has four diagonal blocks $M_{F_i F_i}$, where $F_1 = \{1\}$, $F_2 = \{3\}$, $F_3 = \{6\}$, and $F_4 = \{8, 9\}$, and the linked list is as follows

$$IP = 1, \quad IB = (3, 1, 6, 1, -, 8, 1, 0, 2, -, -)$$

If the index 7 is now added to the set F, then F has three blocks again and the linked list is as follows

$$IP = 1, \quad IB = (3, 1, 6, 1, -, 0, 4, -, -, -, -)$$

As discussed in Section 2, it is necessary in each iteration of the direct methods to solve a linear system with matrix M_{FF} (2 systems for the Parametric Principal Pivoting method) and perform some operations with matrices and vectors. Next, we show that the computational effort of these steps is highly reduced because of the structure of the matrix M.

In Step 0 of the direct methods it is necessary to solve the linear system

$$M_{FF} x = -\mathcal{J} \tag{34}$$

where $F \subseteq \{i : a_i = -\infty \text{ and } b_i = +\infty\}$ and $\mathcal{J} = q_F + M_{FT_1} b_{T_1} + M_{FT_2} a_{T_2}$.

Since M is a tridiagonal matrix, then any component \mathcal{J}_i of the vector \mathcal{J} is given by

$$\mathcal{J}_i = q_i + \sum_{j \in F^1} m_{ij} b_j + \sum_{j \in F^2} m_{ij} a_j$$

where

$$F^1 = T_1 \cap \{i - 1, i + 1\}, \quad F^2 = T_2 \cap \{i - 1, i + 1\}$$

Hence two multiplications are required in the worst case to compute each component \mathcal{J}_i of the vector \mathcal{J} . Since F has the form (32), then the solution of the system (34) reduces to solving the k systems

$$M_{F_i F_i} x = -\mathcal{J}_{F_i}, \quad i = 1, \dots, k$$

Since the matrices of these systems are tridiagonal, then the total number of operations (multiplications and divisions) required to solve the system (34) is

$$3 \sum_{i=1}^k |F_i| - 2k$$

where $|F_i|$ is the number of elements of the set F_i .

We have shown that the computational effort for finding the initial vector \bar{z}_F is greatly reduced. In any other iteration the vector \bar{z}_F can be obtained in an even smaller amount of computational work. To see this, let s be the infeasibility chosen by the direct method. If $s \in F$ (this cannot occur for the Active-Set method) then each one of the direct methods set $F = F - \{s\}$. Let F be given by (32) and $s \in F_r$ for some $r \in \{1, \dots, k\}$. Then the set F_r is replaced by two disjoint subsets F_{r_1} and F_{r_2} (one of these two sets may be empty). Hence the components \bar{z}_i , $i \in F_{r_1} \cup F_{r_2}$ can be obtained by solving at most two linear systems of the form (34) with $F = F_{r_1}$ and $F = F_{r_2}$. On the other hand, as the sets F_i , $i \neq r$, are disjoint with F_r , then all the components of \bar{z}_{F_i} are not modified. So $3 \times (|F_{r_1}| + |F_{r_2}|) - 4$ operations are required to solve such systems. The same can be said for the computation of the vector \bar{p}_F required by the Parametric Principal Pivoting method.

Consider now the case in which the index s chosen by the direct method belongs to T. Since

F is given by (32), then \overline{M}_{F_s} can be computed by solving at most two linear systems with matrices $M_{F_i F_i}$ and $M_{F_{i+1}, F_{i+1}}$, where F_i (F_{i+1}) is not empty if and only if $s - 1 \in F_i$ ($s+1 \in F_{i+1}$). Furthermore

$$\overline{m}_{ss} = m_{ss} + m_{s,s-1} \overline{m}_{s-1,s} + m_{s,s+1} \overline{m}_{s+1,s}$$

that is, it requires at most two multiplications. After this quantity \overline{m}_{ss} is computed, then only the components \overline{z}_j such that $j \in F_i \cup F_{i+1}$ have to be updated by the formulas (12), as the other components of \overline{z}_F do not change their value. In the Parametric Principal Pivoting method, F is increased by the index s. Hence the vector \overline{z}_F (and \overline{p}_F) is computed by just solving a linear system with matrix $M_{F'_1 F'_1}$, where $F'_1 = F_i \cup F_{i+1} \cup \{s\}$.

Finally consider the computation of the vector \overline{w}_T . For any $t \in T$,

$$\overline{w}_t = q_t + \sum_{j \in T_1} m_{tj} b_j + \sum_{j \in T_2} m_{tj} a_j + \sum_{j \in F'} m_{tj} \overline{z}_j$$

where

$$T'_1 = T_1 \cap \{t-1, t, t+1\}$$

$$T'_2 = T_2 \cap \{t-1, t, t+1\}$$

$$F' = F \cap \{t-1, t+1\}$$

This means that each component \overline{w}_t of this vector requires three multiplications. Note that the total computational effort to compute \overline{w}_T can be reduced if the formulas (12) are used.

Since each iteration of each one of the direct methods consists of the computation of the quantities stated above and the updating of the linked list (whenever the set F is modified), then it is clear that the computational effort in each iteration is quite small. This allows the solution of quite large Convex BQPs in a reasonable amount of time even when the number of iterations of the direct method is quite large. This will be confirmed by the computational experience presented in the next section.

The implementation of the Active-Set Method can be used with modifications similar to those stated in the previous section to find stationary points (local minima) of Nonconvex BQPs. Furthermore Concave BQPs do not require the use of the linked list.

5 - Computational Experience and Conclusions

In this section we present a computational investigation of the efficiency of the three direct methods and their implementations discussed in the previous sections. To do this, some Convex BQPs with $M \in \text{SPD}$ have been constructed and are stated below.

- (i) **TP1** - Convex BQP taken from elastoplastic analysis of structures [29].
- (i) **TP2, TP3, TP4, TP5** - Convex BQPs whose matrices are randomly generated by Stewart's technique described in [40].
- (i) **TP6, TP7, TP8** - Convex BQPs which occur in portfolio selection models whose matrices are randomly generated by the technique described in [32]. If n and m are the parameters defined in that paper, then

TP6 - $n = 200, m = 5$

TP7 - $n = 200, m = 20$

TP8 - $n = 600, m = 2$

(iv) **TP9, TP10** - Convex BQPs whose matrices M satisfy $m_{ij} \leq 0$ for all $i \neq j$ and are randomly generated.

(v) **TP11, TP12** - Convex BQPs with pentadiagonal matrices taken from [24].

(vi) **TP13, TP14** - Convex BQPs with tridiagonal randomly generated matrices.

The characteristics of the test problems are described in Table 1. We have decided to present only the results for Convex BQPs with all lower bounds equal to zero, since the performance of the algorithms with infinite or nonzero lower bounds is similar. As far as the upper-bounds is concerned two types of problems have been considered:

(i) all the upper-bounds are infinite (the Convex BQP reduces to a LCP).

(ii) all the upper-bounds are finite and positive.

In both cases the vector q is generated by using a technique described in [36] that fixes the numbers VA and VB of variables that are at the lower and upper bounds in the optimal solution of the Convex BQP. For both the two types of problems the vector $z=0$ is the initial solution ($F=\emptyset$) for the algorithms. Hence VA is the number of variables that are both in the initial and optimal active sets. We also note that for nondegenerate solutions (all the basic z -variables have a value strictly between the bounds) the quantity $n - (VA+VB)$ represents the number of elements of the set F corresponding to the optimal solution of the Convex BQP. If the optimal solution of the Convex BQP is degenerate, then this number is greater than or equal to $n - (VA+VB)$. We have tested problems with nondegenerate (Tables 2,3 and 5) and degenerate optimal solutions (Table 4).

We have tested in the experiments two versions for the Active-Set and Bard-type algorithms. We denote by ASET1 and ASET2 the versions of the Active-Set method that incorporate criteria (15) and (16) respectively. Furthermore the versions of the Bard-type algorithm incorporating criteria (10) and (14) are denoted by BARD1 and BARD2 respectively. We have also found that in general the better choice for the vector p of the Parametric Principal Pivoting Method is as follows

$$\begin{cases} p_i = -1 & \text{if } i \in T_1 \\ p_i = 1 & \text{if } i \in T_2 \\ p_i = 0 & \text{if } i \in F \end{cases}$$

where T_1, T_2 and F are the sets given in Step 0 of this method. All the experiences have been performed on a CDC CYBER 180-830 of the University of Porto, whose machine epsilon [10, chapter 2] is $\epsilon_M = 10^{-14}$. In the tables containing the results of the experiments, there are some parameters with the following meanings:

n = dimension of the Convex BQP = order of the matrix M

nonzer = number of nonzero elements of $M = 2 \times \text{nonz} + n$, where nonz is the number of nonzero elements of M strictly below the diagonal.

stor = total primary storage given by (30).

VA = number of variables at the lower-bound in the optimal solution.

VB = number of variables at the upper-bound in the optimal solution (VB=0 for the LCP).

NO = number of operations (multiplications and divisions) multiplied by 10^{-4} .

NI = number of iterations.

TO = CPU time in seconds for finding an ordering for the rows and columns of M (minimum-degree algorithm).

TF = CPU time in seconds for the symbolic phase.

T = CPU time in seconds for the direct methods.

RES = Residual of the optimal solution. It is given by

$$\| \bar{w} - (q + M\bar{z}) \|_2 = \left[\sum_{i=1}^n \left(\bar{w}_i - q_i - \sum_{j=1}^n m_{ij} \bar{z}_j \right)^2 \right]^{1/2} \quad (35)$$

If α and β are positive integer numbers we use the quantity $\alpha\text{-}\beta$ to represent the Residual $\alpha \times 10^{-\beta}$.

C = Cycling has occurred in the algorithm.

NS = The algorithm has not terminated after 5000 seconds of CPU time but a cycling has not occurred.

Tables 2 and 3 present the results of the experiments on the solution of the test problems by the three direct methods when the optimal solution is nondegenerate. We have tested both the versions BARD1 and BARD2 of the Bard-type method and also the versions ASET1 and ASET2 of the Active-Set algorithm. As expected, the results show that the versions BARD2 and ASET2 usually require less number of iterations than the corresponding versions BARD1 and ASET1. This difference between the number of iterations can even be quite large for some problems. For the LCP whose matrix M has nonpositive off-diagonal elements the number of elements of the set F is always increasing [16] and both the versions 1 and 2 of the direct methods require the same number of iterations. The versions 1 for BQPs with finite upper bounds can even require smaller number of iterations for such matrices.

Despite the number of iterations for the versions 2 being smaller in general, the time T and the number of operations NO do not decrease proportionally. In some cases the versions 1 are more efficient on these parameters (they are the most important) even when the number of iterations is larger. This better performance of the versions 1 can be firstly explained by the search for the infeasibility in each iteration, which obviously requires less work for the versions 1. Furthermore, if the number of deletions that occur in the set F is small (it is zero for the LCPs with $M \in \text{SPD}$ and $m_{ij} \leq 0$) then, by the least-index rule employed by the versions 1, quite often the index s to be added to F in the next iteration satisfies $s > j$ for all $j \in F$. As we have seen, in this case the updating procedure for the LDL^T decomposition has a much lower cost and the number of operations and the resulting CPU time per iteration turn to be much smaller.

Each version of the Active-Set method is in general more efficient, particularly in terms of operations and time, than the corresponding version of the Bard-type method. The reasons for this worse performance of the Bard-type method are presented below.

(i) Since in the Active-Set method the variables z_j of the successive basic solutions

must be kept inside or at the bounds, then deletions in the set F start in the initial stages of the algorithm. In the Bard-type method the set F tends to increase initially and deletions are more frequent in the final stages of the algorithm. Therefore the number of elements of the set F corresponding to the different basic solutions is usually smaller for the Active-Set method. Hence the solution of the linear systems with matrices M_{FF} usually require less computational work for the Active-Set method.

- (ii) For the Active-Set method, whenever a deletion in the set F takes place, only a component of the vector \bar{w}_T has to be computed in the next iteration. This does not occur in the Bard-type method where the whole vector \bar{w}_T is required.
- (iii) In the Active-Set method M_{F_S} is always the right-hand side vector of the system that has to be solved in each iteration. In the Bard-type method the vector q_F is quite often the right-hand side vector. Since M_{F_S} is usually sparser than q_F , then for the same set F the number of operations for solving a system with matrix M_{FF} in the Active-Set method is never greater than that required for solving such system in the Bard-type method.

Another disadvantage of the Bard-type method lies on the possibility of occurrence of cycling for Convex BQPs with some or all finite upper-bounds. These cases are mentioned in Table 3 by a letter C. In general cycling occurs more frequently for the version BARD1.

The Parametric Principal Pivoting (PARM) method share the characteristics of the Bard-type method but has two disadvantages over this method:

- (i) It is necessary to solve two systems with matrix M_{FF} in each iteration.
- (ii) There is always a change in the set F in each iteration. Hence in the solution of Convex BQPs with some finite upper-bounds it is necessary to update the LDL^T decomposition of M_{FF} in each iteration. This may not be required in the other two direct methods, where the infeasibility s may be moved from one of the sets T_i ($i=1, 2$) to the other without modifying F .

These two drawbacks explain the worst performance of the PARM method. However, the algorithm is always convergent to the optimal solution of the Convex BQP when $M \in SPD$.

As stated before, $z=0$ has been chosen as the initial solution for the three direct methods, that is, we have set $T_2 = \{1, \dots, n\}$, $F = T_1 = \emptyset$ initially. We have generated the test problems by introducing the quantities VA and VB representing the number of variables that are at the lower and upper bounds in the optimal solution. Hence $n - VA$ measures the number of variables that move from the initial to the optimal active sets. The three direct methods are quite sensitive to this quantity $n - VA$. In general, the number of iterations NI of the methods BARD2, ASET2 and PARM satisfy

$$n - VA \leq NI \leq \frac{3}{2} (n - VA)$$

and the upper-bound can be $5 (n - VA)$ for the versions ASET1 and BARD1. It is also interesting to mention that $[NI - (n - VA)]$ tends to decrease with an increase of $(n - VA)$.

Table 4 presents the results of the experiments for LCP test problems with degenerate optimal solutions. The results for Convex BQPs with finite upper bounds show a similar behaviour and are not presented. The efficiencies of the algorithms BARD2, ASET2 and PARM do not seem to be much deteriorated with the occurrence of degeneracy in the optimal solution of the Convex BQP. However, the behaviour of the versions ASET1 and BARD1 may be catastrophic in presence of degeneracy in the optimal solution of the Convex BQP. This has actually occurred in the solution of the test problems TP11 and TP12, where no cycling has occurred but the algorithms have not been able to terminate after 5000 seconds of CPU time. Note that for $n = 50$ and the vector q generated by the same technique [36] BARD1 requires 142971 iterations to find the global minimum. It is important to mention that Lin and Pang [24] have reported that the Projected SOR method has not been able to find the global minimum in a reasonable amount of time for a LCP with the same matrix M of order $n = 50$ and a different right-hand side vector q .

To test the efficiency of the implementation for the direct methods with general sparse SPD matrices, we have considered for each test problem the value of three quantities, namely, the total number of operations NO , the CPU time T and the Residual RES defined by (35). To gain a better idea about how good these quantities are for each problem, we have decided to solve for each matrix M a system $Mx + q = 0$, where q is a vector generated by the same technique [36] with $F = \{1, 2, \dots, n\}$. The system has been solved by the General Sparse Solver for SPD matrices described in [11, chapter 5] and the residual of the solution is given by the same formula (35) with $\bar{w} = 0$. The results of the solutions of these systems are presented in Table 1. The comparison of the quantities NO , T and RES for the systems and for the direct methods shows that the implementation is quite efficient in terms of speed and accuracy.

Table 5 contains the results of the experiments of solving Convex BQPs with tridiagonal SPD matrices. Unlike the other tables, only the results for the Active-Set and Bard-type methods are presented, since the PARM algorithm is less efficient. As before we have tested both versions 1 and 2 of these methods. The versions 2 of the algorithms require less number of iterations and operations. However, the CPU time is inferior when the versions 1 are used and this difference increases considerably with the dimension n and the value $(n - VA)$ stated before. The better performance of the versions 1 is a consequence of the block diagonal structure of the matrices M_{FF} . As discussed before, this structure implies that quite few components of the vector \bar{z}_F are modified in each iteration. Since the least index rule is employed by the versions 1 of these direct methods, then the search for the infeasibility s is usually quite small. In fact, this search only starts from the first component of \bar{z}_F whose value has changed in the previous iteration. On the other hand, in the versions 2 the infeasibility s is found by comparing the absolute values of \bar{w}_i (and \bar{z}_i for the Bard-type method). Hence the search is much longer and becomes quite expensive when n is quite large.

As before, each version of the Active-Set method is more efficient than the corresponding version of the Bard-type method. However, the difference is not very big. The reason for this lies on the fact that the number of operations is only related with the subsets of F that are modified in each iteration. Hence the dimensions of the sets F (which are larger for the Bard-type method) have not an effect similar to the case of general sparse matrices. Table 5 also

contains the results of the solution of the linear system $Mz + q = 0$, where q is generated by the technique described in [36] with $F = \{1, \dots, n\}$. As before the comparison of the quantities NO, T and RES for the systems and for the direct methods shows that the implementation is quite good in terms of speed and accuracy. Furthermore the storage space for this implementation is $(3n + 1)$, which corresponds to the diagonal and subdiagonal of M and the linked list.

As a final conclusion of this study we can claim that there is not a big gap between the efficiencies of the Active-Set and Bard-type methods, but the former is consistently more efficient. The versions of these two methods that differ on the choice of the infeasibility in each iteration have a different performance for general sparse and tridiagonal matrices. It seems worthwhile to design partial selection techniques [4,15] which can exploit the benefits of each one of the selection procedures that have been tested in this paper. Cycling may occur in the Bard-type method but the phenomenon is quite rare. The Principal Parametric Pivoting algorithm is convergent but it is in general less efficient than the other two direct methods.

In our experiments we have always used $z = 0$ as the initial solution for all the direct methods. The Active-Set method can start with any vector z satisfying the bounds. The choice of this starting vector has a determinant effect on the efficiency of the Active-Set method, as the computational study presented in this paper has shown. Recently some work has been done on the design of hybrid methods that try to identify the optimal active-set in a reasonable amount of time [26]. The Bard-type and the Parametric Principal Pivoting methods can even start with a vector z which does not satisfy the bounds (the vector p should be chosen with care). This is a great advantage of these methods over the Active-Set method. It seems easier to find a procedure that gets an advanced basic solution that is close to the optimal solution. Recently we have designed an algorithm for solving LCPs with SPD matrices [17] that incorporates such a procedure before using Murty's Bard-type method [28]. This method is in general much more efficient than the direct methods discussed in this paper particularly when the dimension of the LCP is quite large. We have been investigating on the design of a similar procedure to solve Convex BQPs with some or all finite upper-bounds.

It is well-known that there are algorithms for nonlinear optimization, such as Diagonalization and Sequential methods, that rely on the solution of a sequence of Convex BQPs. All the matrices of these Convex BQPs have the same sparsity pattern and in some cases are even equal. Since the Bard-type and Parametric Principal Pivoting methods do not require a vector satisfying the bounds to start with, then they may be much more useful than the Active-Set methods for these types of algorithms. This is another important topic of our research.

An advantage of the Active-Set method over the two remaining algorithms is its ability to solve Convex BQPs with singular SPSD matrices. The Active-Set method can even find a stationary point (local minimum) of a Nonconvex BQP. The extensions of the Bard-type and Parametric Principal Pivoting methods to solve Convex BQPs with singular SPSD matrices and Nonconvex BQPs have not been done unless in some special cases. This is another important topic of future research.

Problem	n	nonzer	stor	ANALYSE		SYSTEM		
				TO	TF	NO	T	RES
TP1	484	9 920	17 898	3.8	0.49	11.1	2.5	2-11
TP2	600	11 442	17 463	3.7	0.21	6.32	1.51	5-12
TP3	800	10 920	16 780	3.2	0.22	4.83	1.28	4-12
TP4	1000	16 386	25 079	5.2	0.32	8.14	2.03	4-12
TP5	1500	9 824	15 468	4.4	0.23	2.99	1.01	3-12
TP6	200	1 922	3 035	0.5	0.05	1.07	0.27	1-12
TP7	200	8 072	16 083	3.2	0.48	27.6	5.27	5-12
TP8	600	9 878	15 135	6.9	0.24	12.9	2.59	3-12
TP9	500	2 156	8 332	2.9	0.33	14.7	2.92	7-12
TP10	1000	2 842	6 217	1.7	0.16	2.59	0.71	6-12
TP11	500	2 494	3 991	---	0.09	0.6	0.26	3-12
TP12	1000	4 994	7 991	---	0.18	1.2	0.51	4-12

TABLE 1 - Matrices of the test problems, ANALYSE Phase and Systems

Problem	N	BARD1			BARD2			ASET1			ASET2			FARM								
		NI	NO	RES	NI	NO	RES	NI	NO	RES	NI	NO	RES	NI	NO	RES						
TP1	484	339	553	116.4	27.6	8-12	391	39.9	10	8-12	509	52.4	15.8	1-11	195	24.7	7.5	1-11	175	67.4	13.2	7-12
		97	637	325.8	74.7	1-11	411	230.4	54.7	1-11	637	122.5	44.7	2-11	611	109.3	37.7	2-11	401	417.9	81.9	2-11
TP2	600	420	554	85.4	22.5	2-12	210	41.3	11	2-12	518	41.6	16.2	2-12	210	29.3	9.1	2-12	242	102.3	20.4	2-12
		120	836	343	88.9	4-12	520	249.9	64.2	4-12	782	80.1	47.3	5-12	538	128.1	48.7	4-12	544	531.1	109.3	4-12
TP3	800	560	602	86	25.8	1-12	278	50.2	15.3	1-12	584	46.3	17.7	2-12	278	36.5	12.9	2-12	294	110.4	24.2	1-12
		160	1004	402.8	112.7	3-12	674	300.9	84.4	4-12	960	96.5	64.3	4-12	694	156.2	65.1	3-12	724	675.5	147.8	4-12
TP4	1000	700	790	178.4	49.6	2-12	334	88.3	25.1	2-12	740	89.2	31.9	2-12	334	66	21.5	2-12	354	198.7	41.1	2-12
		200	1476	837.9	226.3	4-12	838	559.1	149.9	4-12	1188	191.5	123.4	5-12	858	286.3	114.6	4-12	886	1213	256.6	4-12
TP5	1500	1050	920	117.4	190.2	1-12	500	78.2	33.7	1-12	894	69.4	37.4	1-12	502	61.4	31	1-12	526	185.4	51.3	1-12
		300	1600	572.9	209.8	3-12	1250	485.6	174.7	3-12	1564	175.7	136.4	3-12	1260	278.9	147.1	3-12	1306	1086	289.9	3-12
TP6	200	140	108	2.85	1.09	5-13	60	1.88	0.65	7-13	108	2.04	0.88	7-13	60	1.61	0.58	5-13	60	4.19	0.99	6-13
		40	174	7.3	2.7	1-12	160	15.9	4.5	1-12	174	5.3	2.3	2-12	160	10.2	3.3	1-12	160	32.1	7.1	1-12
TP7	200	140	312	113.3	24.9	2-12	60	10.2	2.3	2-12	296	77.7	17.7	3-12	60	9.5	2.2	2-12	60	19.9	3.6	2-12
		40	242	132.5	28.2	4-12	160	128.2	27.9	4-12	242	116.4	25.1	6-12	160	112.5	25.4	5-12	160	210.4	40.2	5-12
TP8	600	420	278	29.5	9.7	1-12	180	30.4	7.8	2-12	278	22	7.6	2-12	180	25.7	6.8	3-12	182	67.5	12.8	2-12
		120	494	66.2	21.9	5-12	480	216.8	51.8	3-12	494	56.2	18.3	6-12	480	135.7	34.9	1-11	480	430.1	82.6	3-12
TP9	500	350	150	4.1	2.33	3-12	150	6.7	3.25	2-12	150	3.6	2.3	2-12	150	5	3	1-11	150	13.7	4.6	2-12
		100	400	41.6	14.7	7-12	400	133.5	35.1	7-12	400	32.5	13.6	2-11	400	92.4	28.9	4-11	400	190.7	45.5	8-12
TP10	1000	700	300	10.1	7.9	2-12	300	15.8	10.5	2-12	300	8.6	7.8	2-12	300	13.3	10	2-12	300	35.6	14.9	2-12
		200	800	72.9	40.8	5-12	800	143.8	61.2	7-12	800	47.6	37.5	6-11	800	74.6	49.9	3-10	800	277	90.8	6-12
TP11	500	350	386	14.9	6.9	1-12	150	5.7	3.09	1-12	384	8.9	5.2	1-12	150	5.1	3.08	1-11	150	13	4.5	1-12
		100	714	71.2	28.1	3-12	400	41.6	20	3-12	566	21.9	16.6	7-12	400	26.9	18.2	1-11	404	88.1	30.1	3-12
TP12	1000	700	802	61.6	27.9	1-12	300	22.6	12.1	1-12	792	35.9	20.7	2-12	300	20	12.1	1-12	300	51.9	17.6	1-12
		200	1488	275.3	107.9	5-12	802	163.5	77.5	4-12	1118	79.6	63.9	1-11	802	102.9	69.8	5-12	816	355.3	120.8	4-12

TABLE 2 - LCPs with nondegenerate solution

N	V/A	V/E	EARD1			EARD2			ASET1			ASET2			FARM								
			NI	NO	I	RES	NI	NO	I	RES	NI	NO	I	RES	NI	NO	I	RES					
TP1	484	340	72	752	79.6	21.8	5-12	214	19.7	6.9	6-12	709	34.9	11.9	1-11	224	10.4	4.2	9-12	251	103.3	22.2	6-12
			96	194	1176	294.8	65.7	1-11	556	150.	49.9	1-11	1115	120.9	40.2	2-11	580	67.4	23.5	2-11	617	725.1	150.7
TP2	690	420	90	901	61.3	21.6	1-12	277	17.3	7.9	1-12	809	30.6	12.5	2-12	300	13.6	6.2	1-12	347	152.	36.5	1-12
			120	240	C			688	92.2	31.9	5-12	1448	107.3	42.7	3-12	807	85.1	27.9	3-12	925	884.7	189.7	4-12
TP3	890	560	120	841	40.8	18.5	1-12	326	16.1	9.7	1-12	797	24.4	11.9	1-12	340	12.1	7.1	1-12	409	162.9	39.5	1-12
			160	320	1780	247.2	81.4	3-12	920	117.4	46.9	3-12	1585	100.4	48.2	3-12	990	92.2	36.	3-12	1149	1055.	239.1
TP4	1090	700	150	1199	105.9	41.	1-12	C			1083	54.2	23.4	1-12	433	25.7	12.5	1-12	536	323.5	73.9	1-12	
			200	400	2400	491.7	151.8	3-12	1102	198.7	74.9	3-12	2181	212.2	92.9	3-12	1233	175.3	64.1	3-12	1472	2021.	446.7
TP5	1500	1050	225	C			588	24.3	25.9	8-12	1127	33.9	25.6	1-12	588	17.8	19.1	1-12	791	295.	92.9	8-12	
			300	600	C			1595	175.4	115.7	2-12	2188	130.1	96.9	2-12	1655	129.6	90.9	2-12	2119	1817.	495.2	2-12
TP6	200	140	30	135	2.56	1.14	6-13	69	0.58	0.47	5-13	133	1.42	0.72	7-13	71	0.39	0.34	6-13	91	7.4	2.	7-13
			40	80	246	7.3	2.8	8-13	169	3.98	1.85	1-12	237	3.96	1.89	1-12	177	1.83	1.22	1-12	241	53.8	13.3
TP7	200	140	30	453	117.4	26.9	2-12	79	9.3	2.4	1-12	450	66.8	16.1	2-12	86	4.81	1.3	2-12	93	36.3	7.3	2-12
			40	80	636	283.8	62.9	4-12	194	82.8	18.8	4-12	645	199.4	45.2	5-12	236	38.6	8.8	4-12	243	367.2	74.3
TP8	600	420	90	332	38.1	12.	2-12	188	2.26	2.98	2-12	311	16.8	6.5	4-12	191	2.08	2.39	1-12	277	110.1	25.5	1-12
			120	240	628	82.2	24.4	2-12	488	23.8	12.3	3-12	629	46.1	16.4	7-12	495	15.5	9.4	5-12	721	746.8	167.5
TP9	500	350	75	182	2.6	2.22	2-12	182	3.6	3.5	2-12	180	1.5	1.6	1-12	181	1.7	2.27	2-12	228	24.1	9.	2-12
			100	200	579	72.3	24.2	3-12	548	170.7	44.9	3-12	574	43.4	17.7	5-11	548	81.8	27.5	3-10	601	380.5	92.9
TP10	1000	700	150	344	4.9	6.6	1-12	345	6.4	10.6	1-12	342	3.3	5.	1-12	344	3.6	7.3	1-12	451	61.8	20.5	1-12
			200	400	1066	66.4	43.4	3-12	1039	97.9	64.6	4-12	1059	26.7	32.1	5-12	1029	32.3	39.1	2-10	1201	488.4	124.3
TP11	500	350	75	438	7.2	5.2	8-13	175	2.5	2.9	7-13	399	1.7	3.2	7-13	174	1.6	2.1	7-13	228	22.9	9.	7-13
			100	200	951	39.	19.6	2-12	589	38.9	23.2	2-12	667	13.1	10.2	2-12	552	14.8	12.9	2-12	633	160.	54.4
TP12	1000	700	150	883	28.9	21.4	1-12	355	9.9	11.7	8-13	819	14.5	12.9	1-12	350	5.9	8.2	1-12	465	97.5	36.1	8-13
			200	400	1795	151.9	75.	2-12	1171	147.1	88.8	2-12	1264	53.9	40.6	3-12	1126	55.7	50.8	3-12	1262	616.5	212.5

TABLE 3 - BOPs with finite upper-bounds and nondegenerate optimal solution

Problem	N	BARSD1			BARSD2			ASES1			ASES2			PARM								
		NI	NO	I	RES	NI	NO	I	RES	NI	NO	I	RES	NI	NO	I	RES					
TP1	484	339	54.8	124.7	29.5	3-11	200	44.3	11.	1-11	524	59.9	17.6	1-11	222	31.4	9.3	2-11	170	64.4	12.5	2-10
		97	657	360.7	81.9	2-11	434	233.9	54.5	2-11	671	133.6	49.	2-11	422	114.2	39.9	2-11	401	422.8	81.9	2-11
TP2	600	420	586	96.9	25.3	3-12	212	41.1	10.9	2-12	537	44.4	15.6	3-12	213	30.1	9.1	2-12	243	105.9	20.6	4-12
		120	902	389.5	101.	7-12	554	284.7	71.9	1-11	877	84.4	52.6	6-12	570	135.6	51.8	6-12	556	560.5	112.	6-12
TP3	800	560	608	89.7	26.5	2-12	281	51.1	15.2	2-12	583	46.3	18.	2-12	281	37.5	13.	2-12	305	117.6	25.2	1-12
		160	1077	444.8	122.	7-12	699	324.7	88.7	5-12	1050	97.7	70.	4-12	717	162.6	67.8	4-12	755	732.6	157.	5-12
TP4	1000	700	843	203.4	55.4	4-12	335	88.9	24.7	3-12	781	94.	34.6	3-12	337	67.9	21.7	3-12	360	204.2	41.7	3-12
		200	1539	905.7	238.	1-11	875	613.5	160.	6-12	1461	108.9	128.	6-12	899	298.6	119.	7-12	933	1344.	278.	5-12
TP5	1500	1050	943	125.8	51.8	3-12	507	79.9	33.6	2-12	907	72.4	38.5	2-12	506	63.5	31.	2-12	530	189.7	52.6	2-12
		300	1656	575.2	200.	1-11	1284	516.4	179.	7-12	1627	174.1	140.	5-12	1301	288.9	150.1	3-12	1328	1134.	292.	5-12
TP6	200	140	97	1.96	0.28	7-13	60	1.88	0.6	7-13	97	1.44	0.65	5-13	60	1.61	0.54	5-13	60	4.21	0.94	6-13
		40	184	10.5	3.4	1-12	160	15.9	4.2	1-12	194	7.54	2.8	1-12	160	10.1	3.1	1-12	160	32.2	6.7	1-12
TP7	200	140	329	126.3	26.1	2-12	60	10.2	2.2	2-12	298	80.1	17.1	2-12	60	9.5	2.1	2-12	61	20.8	3.5	1-12
		40	240	137.8	27.8	5-12	160	128.2	26.3	4-12	240	125.3	25.6	6-12	160	112.5	23.9	5-12	160	210.4	38.2	5-12
TP8	600	420	301	36.9	10.9	3-12	180	30.4	7.6	2-12	288	24.2	7.9	3-12	180	25.7	6.6	3-12	182	67.5	12.4	2-12
		120	539	90.5	27.2	4-12	480	216.8	50.5	3-12	542	78.6	23.9	1-11	480	125.7	34.1	1-11	480	430.3	80.9	4-12
TP9	500	350	150	4.1	2.2	3-12	150	6.7	3.1	2-12	150	3.6	2.2	2-12	150	5.	2.9	1-11	150	13.7	4.3	2-12
		100	400	41.6	13.9	7-12	400	133.5	33.3	7-12	400	32.5	12.8	2-11	400	92.4	27.3	4-11	400	190.7	43.3	8-12
TP10	1000	700	300	10.1	7.4	2-12	300	15.8	9.9	2-12	300	8.6	7.3	2-12	300	13.3	9.5	2-12	300	35.6	14.	2-12
		200	800	72.9	38.1	5-12	800	143.8	57.9	7-12	800	47.6	35.2	5-11	800	74.6	46.6	3-10	800	270.	85.4	6-12
TP11	500	350	————	N S	————	151	5.8	2.9	1-12	1484	86.7	25.9	1-11	151	5.2	2.8	1-12	151	13.2	4.3	1-12	
		100	————	N S	————	411	45.1	22.4	6-12	————	————	N S	————	————	————	428	33.4	23.7	8-12	415	94.	32.9
TP12	1000	700	————	N S	————	300	22.6	11.3	1-12	3597	345.	117.2	4-11	300	20.	11.3	1-12	301	52.3	16.9	1-11	
		200	————	N S	————	825	175.8	83.9	7-12	————	————	N S	————	————	————	849	199.5	87.6	2-11	843	385.3	141.

TABLE 4 - LCPs with degenerate solution

Problem	N	SYSTEM				VA		VB		BARD1			BARD2			ASET1			ASET2				
		NO	T	RES		NI	NO	T	RES	NI	NO	T	RES	NI	NO	T	RES	NI	NO	T	RES		
TP13	1000	0.49	0.08	1-13	LCP	700	---	414	0.2	0.41	4-13	299	0.12	2.4	4-13	407	0.17	0.39	4-13	299	0.1	2.1	4-13
						200	---	878	1.58	0.89	9-13	799	0.93	6.5	9-13	871	1.19	0.84	1-12	802	0.71	5.6	9-13
					BQP	700	150	436	0.24	0.25	4-13	319	0.16	4.4	4-13	436	0.2	0.24	4-13	319	0.14	3.8	4-13
						200	400	1047	0.88	0.89	7-13	929	0.79	12.9	7-13	1039	0.72	0.85	8-13	932	0.57	10.	7-13
TP14	5000	2.49	0.41	3-13	LCP	3500	---	2174	1.06	8.2	8-13	1501	0.59	57.1	8-13	2160	0.94	8.	9-13	1501	0.41	49.	8-13
						1000	---	4456	7.84	13.2	2-12	4007	4.79	156.6	2-12	4442	6.54	12.9	2-12	4009	3.92	120.	2-12
					BQP	3500	750	2290	1.24	3.62	7-13	1587	0.77	106.3	7-13	2281	1.08	3.56	8-13	1587	0.68	92.5	7-13
1000	2000	5486	4.7	15.3		2-12	4660	3.9	311.8	2-12	5671	3.8	15.	2-12	4659	2.8	247.	2-12					

TABLE 5 - LCP and BQP with finite upper-bounds and tridiagonal matrices

References

- [1] Bennett, J. M., 1965, *Triangular factors of modified matrices*, Numerische Mathematik 7, 217-221.
- [2] Cottle, R.W. and Goheen M.S., 1978, *A special class of large quadratic programs*, in "Nonlinear Programming 3", edited by O.L. Mangasarian, R.R. Meyer and S.M. Robinson, Academic Press, New York, 361-390.
- [3] Crabtree, D.E. and Haynsworth, E.V., 1969, *An identity for the Schur Complement of a matrix*, Proceedings of American Mathematical Society 22, 364-366.
- [4] Crowder, H. and Hattingh, J. M., 1975, *Partially normalized pivot selection in linear programming*, Mathematical Programming Study, 4, 12-25 .
- [5] Cryer, C.W., 1983, *The efficient solution of linear complementarity problems for tridiagonal Minkowski matrices*, ACM Transactions on Mathematical Software 9, 199-214.
- [6] Dembo, R. S. and Tulowitzki, U., 1983, *On the minimization of a quadratic function subject to box constraints*, Working Paper Series B# 71, School of Organization and Management, Yale University.
- [7] Duff, I.S., 1984, *Direct methods for solving sparse systems of linear equations*, SIAM Journal Scientific Statistical Computing 5, 605-619.
- [8] Fletcher, R. and Jackson, M.P, 1974, *Minimization of a quadratic function subject only to upper and lower bounds*, Journal Institute Mathematics and Applications 14, 159-174.
- [9] Fletcher, R., and Powell, M.J.D., 1974, *On the modification of the LDL^T factorizations*, Mathematics of Computation 28, 1067-1087.
- [10] Forsythe, G.E., Malcolm, M.A. and Moler, C.B., 1977, *Computer methods for mathematical computations*, Prentice-Hall, Englewood Cliffs, New Jersey.
- [11] George, A. and Liu, J.W.H., 1981, *Computer solution of large positive definite systems*, Prentice-Hall, Englewood Cliffs, New Jersey.
- [12] Gill, P.E., Murray, W. and Wright, M.H., 1981, *Practical Optimization*, Academic Press, New York.
- [13] Glowinski, R., 1978, *Finite elements and variational inequalities*, MRC Technical Report 1885, Mathematics Research Center, University of Wisconsin - Madison.
- [14] Graves, R.L., 1967, *A principal pivoting simplex algorithm for linear and quadratic programming*, Operations Research 15, 482-494.

- [15] Harris, P.M.J., 1975, *Pivot selections methods of the Devex LP code*, Mathematical Programming Study 4, 30-57.
- [16] Júdice, J. J. and Pires, F.M., 1990, *A Bard-type method for a generalized linear complementarity problem with a nonsingular M-matrix*, Naval Research Logistics, 37, 279-297.
- [17] Júdice, J. J. and Pires, F.M., 1988/89, *Bard-type methods for the linear complementarity problem with symmetric positive definite matrices*, IMA Journal of Mathematics Applied in Business and Industry 2, 51-68.
- [18] Júdice, J. J. and Pires, F.M., 1986, *Direct methods for the solution of linear complementarity problems with symmetric positive semi-definite matrices*, Investigação Operacional 6, 115-152.
- [19] Keller, E.L., 1973, *The general quadratic optimization problem*, Mathematical Programming 5, 311-337.
- [20] Kuhn, H. and Tucker, W., 1951, *Nonlinear programming*, in "Second Berkeley Symposium in Mathematical Statistics and Probability", edited by J.Neyman, University of California Press, California, 80-90.
- [21] Law, K.H., 1985, *Sparse matrix modification in structural reanalysis*, International Journal for Numerical Methods in Engineering, 21, 37-63.
- [22] Lawson, C.L. and Honson, R.J., 1974, *Solving least squares problems*, Prentice-Hall, Englewood Cliffs, New Jersey.
- [23] Lin, Y.Y. and Cryer, C.W., 1985, *An alternating direction implicit algorithm for the solution of linear complementarity problems arising from free boundary problems*, Applied Mathematics and Optimization, 13, 1-17.
- [24] Lin, Y.Y. and Pang, J.S., 1987, *Iterative methods for large convex quadratic programs: a survey*, SIAM Journal on Control and Optimization, 25, 383-411.
- [25] Lötsdet, P., 1984, *Solving the minimal least squares problem subject to bounds on the variables*, BIT 24, 206-224.
- [26] Moré, J.J. and Toraldo, G., 1989, *Algorithms for bound constrained quadratic programming problems*, Numerische Mathematik 55, 377-400.
- [27] Murty, K.G. 1983, *Linear Programming*, John Wiley & Sons, New York.
- [28] Murty, K.G., 1974, *Note on a Bard-type scheme for solving the complementary problem*, Opsearch 11, 123-130.
- [29] Neves, A.S., Faustino, A.M., Pires, F.M. and Júdice, J.J., 1986, *Elastoplastic analysis of structures and linear complementarity*, in "OR Models on

- Microcomputers", edited by Coelho, J.D. and Tavares, L.V., North-Holland, Amsterdam, 217-228.
- [30] O'Leary, D.P., 1980, *A generalized conjugate gradient algorithm for solving a class of quadratic programming problems*, Linear Algebra and Applications 34 , 371-399.
- [31] Osterby, O. and Zlatev, Z., 1983, *Direct methods for sparse matrices*, Lecture Notes in Computer Science, 157, Springer-Verlag, Berlin.
- [32] Pang, J.S., 1980, *A new and efficient algorithm for a class of portfolio selection problems*, Operations Research, 28, 754-767.
- [33] Pang, J.S., Kaneko, I. and Hallman, W.P., 1979, *On the solution of some (parametric) linear complementarity problems with applications to portfolio analysis, structural engineering and graduation*, Mathematical Programming, 16, 325-247.
- [34] Pang, J.S. and Lee, S.C., 1981, *A parametric linear complementarity technique for the computation of equilibrium in a single commodity spatial model*, Mathematical Programming, 20, 81-102.
- [35] Plemmons, R.J., 1977, *M-matrix characterizations: I - nonsingular M-matrices*, Linear Algebra and its Applications, 18, 175-188.
- [36] Ramarao, B. and Shetty, C.M., 1984, *Application of disjunctive programming to the linear complementarity problem*, Naval Research Logistics Quarterly, 31, 589-600.
- [37] Sargent, R.W.H., 1978, *An efficient implementation of the Lemke algorithm and its extension to deal with upper and lower bounds*, Mathematical Programming Study, 7, 36-54.
- [38] Van der Laan, G. and Talman, A.J.J., 1989, *An algorithm for the linear complementarity problem with upper and lower bounds on the variables*, Journal of Optimization Theory and Applications, 62, 151-163.
- [39] Westbrook, D.R., 1982, *Contact problems for the elastic beam*, Computers and Structures, 15, 473-479.
- [40] Yang, E.K. and Tolle, J.W., 1985, *A class of methods for solving large convex quadratic programs subject to box constraints*, Working Paper, Management Sciences Department, University of Massachusetts at Boston.

CICLOS E TENDÊNCIA EM SÉRIES ECONÓMICAS: O PIB PORTUGUÊS DE 1913 A 1986

João A.O. Soares
J.A. Assis Lopes
Secção Autónoma de Economia e Gestão
Instituto Superior Técnico
Universidade Técnica de Lisboa

Resumo: É discutido neste artigo o problema da mensuração dos ciclos económicos, e nomeadamente dos ciclos de crescimento, enquanto dependente da caracterização da tendência de longo prazo como determinística ou estocástica - modelos 'TS' ('Trend-Stationary') ou 'DS' ('Difference-Stationary'). Referem-se os instrumentos estatísticos associados a ambos - ajuste de uma função determinística do tempo e análise espectral versus modelização ARIMA ou estrutural.

Os mesmos são utilizados para exame do PIB português 1913-86, concluindo-se pela caracterização 'DS'.

Abstract: This paper discusses the problem of measuring business cycles, namely the growth cycles, as depending on the characterization of the trend as deterministic or stochastic - 'TS' ('Trend-Stationary') or 'DS' ('Difference-Stationary'). Statistical tools associated with both are mentioned - deterministic function of time and spectral analysis versus ARIMA or structural models. Tools are used to examine the portuguese GDP 1913-86, the conclusion being to prefer the 'DS' characterization.

1 - Ciclicidade e Tendência em Séries Económicas

O problema da existência e mensuração dos ciclos económicos ou ciclos de negócios ('business cycles'), usando a tradução literal da designação encontrada na literatura anglo-saxónica, é um problema com mais de cem anos na Teoria Económica, tendo dado origem a vasta literatura teórica e empírica.

Querendo encontrar uma definição de ciclo económico amplamente difundida e suficientemente abrangente para merecer um amplo consenso, pode-se referenciar Arthur Burns e Wesley Mitchell ("Measuring Business Cycles", National Bureau of Economic Research, N.York, 1947). Estes, sinteticamente, caracterizam os ciclos como flutuações da actividade económica, compostos por expansões e recessões que se sucedem ininterruptamente, de duração entre 1 e 10 a 12 anos, sem periodicidade definida.

Esta definição, por sua vez, fica incompleta se não restringermos entre o que se convencionou designar por 'ciclos clássicos' e os 'ciclos de crescimento'.

Os primeiros correspondem à delimitação das fases de expansão e contracção a verificar nas séries brutas representativas da actividade económica. Do ponto de vista estatístico, esse trabalho não oferece particular dificuldade, havendo só que, cuidadosamente, distinguir pequenas flutuações conjunturais, de fases do ciclo propriamente ditas.

Quanto aos 'ciclos de crescimento', a sua referência aparece nos anos 60, quando a economia dos principais países desenvolvidos conhece uma grande fase de desenvolvimento, parecendo então afastadas as contracções da actividade económica, enquanto diminuição absoluta do produto real e aumento do desemprego dos factores produtivos. Em contrapartida, registam-se fases de atenuação das taxas de crescimento, desacelerações da actividade, que não recessões no sentido actual. Os 'ciclos de crescimento' pretendem então exprimir fases de desvio, flutuações, face à tendência de longo prazo. Implicitamente daquelas componentes nas séries económicas.

Ora, isolar, determinar, a componente cíclica manifesta nas séries relevantes da actividade económica, nomeadamente no Produto, implica a assunção de uma decisão quanto à modelização da tendência - assume esta um carácter determinístico ou estocástico.

A opção pela modelização determinística da tendência traduz-se no ajustamento de uma função do tempo - tipicamente a exponencial ou um polinómio de grau baixo-, procedendo-se de seguida à substracção da mesma relativamente à série original. A ciclicidade dos resíduos é então objecto de análise através do uso do periodograma ou, pondo a tónica no seu carácter estocástico, através do ajustamento de um modelo ARMA, nomeadamente um AR(2) com raízes complexas.

Mas esta caracterização das séries económicas, que é consagrada na literatura com a designação de modelos 'Trend-Stationary' (TS), sofre vasta contestação. Pode-se citar, a propósito, o artigo de Nelson e Kang (1981) que mostra como o facto de se retirar uma tendência determinística a um processo tipo passeio aleatório pode introduzir falsa periodicidade nos correspondentes resíduos. Também noutro artigo do mesmo ano, Beveridge e Nelson propõem mais explicitamente uma forma de decomposição de séries económicas baseada em modelos ARIMA. Finalmente, Nelson e Plosser (1982) consagram a designação de 'DIFFERENCE-STATIONARY' (DS) para processos em que as 1^{as} (ou outras) diferenças da série original ou do seu logaritmo, são um processo estacionário modelizável por um modelo ARMA.

Nesse mesmo artigo, os autores analisam diversas séries da economia norte-americana, com destaque para a do PIB, e concluem favoravelmente pela sua caracterização como enquadrando-se nos modelos 'DS', i.e., pela tendência estocástica.

Harvey (1984), por sua vez, embora criticando alguns aspectos do artigo acima referido, chega a conclusões concordantes com este, socorrendo-se da modelização estrutural em espaço de estados.

No presente trabalho utilizar-se-ão as diferentes ferramentas estatísticas referidas para análise do caso da economia portuguesa.

A série utilizada - PIB de 1913 a 1986, a preços de 1977 - resulta da compatibilização das estimativas referentes ao período 1913-47 (N. Valério, 1983), com os dados de 1948-57 do INE (INE, 1959), de 1958-1978 do Banco de Portugal (B. Portugal, Doc. Trabalho, nº15), e por fim com aqueles publicados pelo INE de 1979 até 1986. Uma explicação mais detachada sobre a mesma pode ser encontrada em Soares (1989).

2 - Análise Estatística do PIB Português 1913-86

2.1. - Análise Gráfica e dos Valores Brutos

Comece-se a análise pelo gráfico da série - Fig.1.

Ele revela essencialmente 3 fases:

- Até finais dos anos 20: Zona de declive quase nulo, com marcada flutuação conjuntural;
- até 1974: Zona de crescimento, primeiro mais lento e depois mais vivo (a partir sobretudo do fim da 2ª Guerra Mundial). Aqui os únicos anos de decréscimo do Produto real são: 1936, 1942, 1944 e 1945.
- de 1974 a 1986: novamente maior pronunciamento cíclico acompanhada de uma tendência crescente. Dois valores acentuados: 1975 e 1983/84.

Nesta evolução, que aliás apresenta similitudes com a da economia internacional, é possível estabelecerem-se os picos e vales delimitadores das fases de expansão e recessão, sem levar em conta os desvios à tendência - ciclos clássicos. Podem-se, então, calcular os

Intervalo médio entre máximos: 7.44 anos (D.P. = 9.11 anos)

Intervalo médio entre mínimos: 7.33 anos (D.P. = 8.8 anos)

As medidas dos intervalos entre extremos, consonantes com as diferenças das fases já referidas e em que sobressai um período de crescimento contínuo de 1945 a 1974, não indicam a existência de regularidade estatística para todo o período em análise.

PIB 1913 - 1986

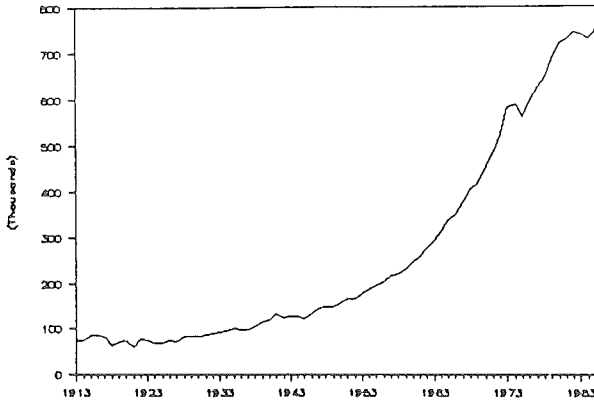


Fig. 1

2.2. - Ajustamentos de uma Tendência Determinística

Procedendo de seguida ao ajustamento de uma tendência determinística, ensaiou-se a utilização de exponencial e de polinómios de grau 2 a 4. Aquela fornece um valor de \bar{R}^2 (R^2 ajustado) de 0.952, enquanto para este obtiveram-se valores de 0.985, 0.988 e 0.991. Refira-se, relativamente ao ajuste da exponencial, que o valor obtido do parâmetro da tendência foi de 3.7. Este pode ser interpretado, como é sabido, com a taxa anual de crescimento percentual de tendência, ao longo de todo o período.

A primeira conclusão que se pode extrair dos valores é que, embora se saiba que para séries económicas suficientemente longas a tendência de longo prazo é largamente dominante, é de sublinhar que o ajustamento polinomial deixa por explicar somente cerca de 1% da variância total - isto para a componente cíclica mais a irregular. Tal valor pode ser contraposto os 5% da autoregressivos de 1ª ordem. A análise do periodograma confirmou o domínio referente a baixas frequências.

- A segunda questão que faltava focar prende-se com a análise directa dos resíduos dos ajustamentos. Embora ela não distinga a componente cíclica da irregular, corresponde à concepção lata de flutuações do Produto, interessante na comparação com os resultados encontrados para a série bruta. Aí é de notar, sobretudo, as novas divisões introduzidas pelo corte efectuado pela linha de tendência à grande fase intermédia de crescimento - 1945 a 1974. O resultado é a diminuição do intervalo médio entre máximos relativos que passa para valores próximos de 4 anos, com D.P. de cerca de 1.7 anos, para a generalidade dos ajustamentos.

2.3. - Tendência Estocástica - Modelos ARIMA

Neste ponto procedeu-se à logaritmização e diferenciação de série, por forma a torná-la estacionária quanto às suas média e variância. Na Fig.2 observa-se que se corrigiu no essencial o problema da estacionaridade em torno da média, mas é nítida a maior variância na fase inicial.

As funções estimadas de autocorrelação amostral (V. Quadro 1), apresentam, por seu turno, uma estrutura oscilante de valores positivos e negativos, não significativos a 2 D.P., a qual aponta para um processo aleatório, exceptuando-se ligeiramente as 1ª e 3ª autocorrelações. A primeira assume um valor negativo ao contrário da 3ª. Este comportamento decorre do carácter muito oscilante da série diferenciada, onde se denota uma ciclicidade muito curta, com distância média entre picos de cerca de 3 anos.

Avançando na análise, entendemos proceder ao seccionamento da série no ano de 1925, correspondendo às fases distintas que tínhamos detectado. A primeira secção fica com um número de pontos demasiado pequeno para prosseguir o estudo. A fase mais recente vem revelar um padrão de ruído branco. Tal se evidencia no correlograma, correlograma parcial e ainda no periodograma integrado que se calculou.

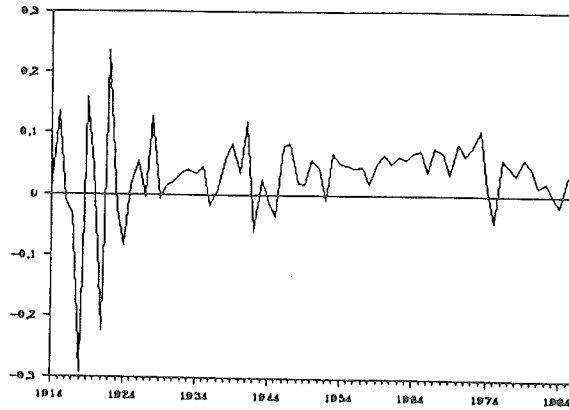
$\Delta \text{LOG} (\text{PIB } 1913 - 86)$ 

Fig. 2

Este último, embora já na série completa ficasse dentro das bandas de 95%, correspondentes ao teste de Kolmogorov - Smirnov, é agora mais nitidamente uniforme.

Por sua vez refira-se que o \bar{R}^2 do $\log(\text{PIB})$ sobe, no período de 1925-1985, para 0.984.

A conclusão que se extrai daqui para a presença de um modelo 'DS' no caso, um passeio aleatório do $\log(\text{PIB})$. Tal merece ser confrontado com os modelos estruturais de Harvey, exponencial. Mas neste caso, a mera análise gráfica parece demonstrar ser irrealista a assunção de taxa de crescimento anual constante para todo o período que lhe é implícita.

Passemos agora à análise dos resíduos:

A visão geral da série residual pode ser obtida, no que concerne à exponencial, se referirmos que a curva ajustada vai cortar a série original em meados dos anos 20 e 60. Marca então 3 grandes fases, acima, abaixo, e novamente acima da tendência. Quanto aos ajustamentos polinomiais, a aderência nas extremidades é superior (daí o maior \bar{R}^2), verificando-se ainda a elevação do pico de 1941 e anos adjacentes.

Com este padrão é óbvio que os periodogramas dos resíduos, enquanto estimadores do espectro teórico, são largamente dominados por ciclos correspondentes a 1 ciclo para a exponencial e 2 ciclos para os ajustamentos polinomiais - períodos de 74 e 37 anos. Para se dar uma ideia mais precisa refira-se que, para a exponencial, esse pico corresponde a cerca de 74% da variância total, sendo p.ex., de 42% para o polinómio de 3º grau. Na generalidade não deparámos com mais nenhum pico correspondendo a uma explicação de variância superior a 10%.

Mas qual é a principal ideia a extrair da análise?

Pensamos que é a existência de um pico no periodograma para baixas frequências, correspondendo a valores para o período entre 100 e 50% da dimensão da série. Tal parece exemplificar o "característico espectro das séries económicas", correspondente ao que Nelson e Kang apontaram como resultado de inapropriada extração de tendência. No artigo já referido, de 1981, eles mostram como uma passeio aleatório a que se retire uma tendência determinística, apresenta uma função de densidade espectral com um pico dominante correspondente a um período de cerca de 80% da dimensão da amostra.

Interessante também é comparar os valores das 1^{as} seis autocorrelações dos resíduos referentes à regressão do logaritmo do PIB (i.e., resíduos do ajuste da exponencial), de um passeio aleatório (dimensão = 100) e do PNB real norte-americano entre 1909 e 1970. As duas últimas séries foram extraídas de Nelson e Plosser (1982, pág.150, tabela 4).

Por fim, quanto à análise dos resíduos resultantes da extração de uma tendência determinística, duas ordens de questões nos merecem reparo:

- A primeira refere-se à possibilidade de ajustar funções polinomiais ao logaritmo do PIB, e não directamente a este. É a via prosseguida por Pinto Barbosa (1985) numa análise do PIB, entre 1953 e 1980, tendo por fonte o FMI.

Quadro 1
Autocorrelação Amostras dos Resíduos

Série	r1	r2	r3	r4	r5	r6
PIB português	0,88	0,81	0,73	0,61	0,52	0,49
passeio aleatório	0,91	0,82	0,74	0,66	0,58	0,51
PNB EUA	0,87	0,66	0,46	0,26	0,19	0,07

O autor ensaia o ajustamento de uma recta e de um polinómio de 2º grau para o referido log (PIB). Ambos os ajustamentos apresentam R^2 de 0,991, sendo escolhido o modelo menos parcimonioso para ajustamento de um modelo ARIMA, apesar do reduzido número de observações. É então identificado um AR(2), indicativo de ciclicidade estocástica, com período de 4,16 anos.

No entanto, com os parâmetros apresentados e a fórmula para o período 'd' corrigida (V.Box e Jenkins, Holden Day, ed. 1970, pág. 60 ou Gottman, C.U.P., 1981, pág. 233),

$$d = \frac{2\pi}{\arcsin [|\theta_1|/(2(-\theta_2)^{1/2})]}$$

O valor obtido foi de 17,34 anos.

De qualquer maneira, o que é interessante referir é que não parece existir justificação teórica para o ajuste do logaritmo que não da série original. Do ponto de vista prático tentou-se comprovar essa ideia usando os nossos dados. Os R^2 encontrados foram muito idênticos aos dos ajustamento directo e a estrutura dos correlogramas e correlogramas parciais mostrou-se também semelhante, embora com valores mais atenuados apontando para processos

Quadro 2
 $\Delta \text{LOG}(\text{PIB } 1913-86)$ $\Delta \text{LOG}(\text{PIB } 1925-86)$

Desfa- samento	$\Delta \text{LOG}(\text{PIB } 1913-86)$				$\Delta \text{LOG}(\text{PIB } 1925-86)$			
	Autoc. Amost.	DP	Autocor. Parciais Amostrais	DP	Autoc. Amostrais	DP	Autoc. Parciais Amostrais	DP
1	-.22701	.11704	-.22701	.11704	-.00014	.12804	-.00014	.12804
2	-.11174	.12292	-.17215	.11704	.02026	.12804	.02026	.12804
3	.31259	.12431	.26688	.11704	.06462	.12809	.06465	.12804
4	-.12653	.13465	-.01356	.11704	-.14030	.12862	-.14126	.12804
5	-.04496	.13627	-.02008	.11704	-.06349	.13111	-.06687	.12804
6	.17441	.13647	.08273	.11704	.22058	.13161	.22945	.12804
7	.08920	.13949	.19919	.11704	.13875	.13754	.17096	.12804
8	-.00239	.14027	.10308	.11704	.05619	.13981	.03066	.12804
9	.07786	.14027	.05963	.11704	.12133	.14018	.06665	.12804
10	-.08999	.14086	-.13295	.11704	-.06165	.14189	-.02539	.12804
11	.00828	.14164	-.02370	.11704	-.12771	.14233	-.08168	.12804
12	.00793	.14165	-.04653	.11704	-.05232	.14420	-.09126	.12804
13	.11643	.14166	.15941	.11704	.13794	.14451	.13347	.12804
14	-.08033	.14296	-.09302	.11704	-.13566	.14665	-.17073	.12804
15	.01501	.14358	-.02542	.11704	.04137	.14870	-.06226	.12804

2.4. - Modelos Estruturais

Procedeu-se ao ajustamento de 3 modelos para o logaritmo do produto:

- tendência estocástica (nível + crescimento + componente irregular);
- tendência cíclica (tendência estocástica com ciclo incorporado + comp.irregular);
- tendência + ciclo (tendência estocástica + ciclo estocástico + componente irregular).

Os modelos a) e b) fornecem ambos o valor de $R^2 = 0,9937$, $R_D^2 = 0,09$ (i.e. 9 pontos percentuais acima do passeio aleatório + taxa de impulso) e os correlogramas dos resíduos apontam para ruído branco. Por sua vez as variâncias associadas à tendência, muito próximas de zero, parecem reflectir a grande supremacia da persistência na evolução da série, considerada neste período longo.

No que concerne ao modelo b) é de salientar também que o período estimado de ciclicidade é de 6,57 anos, sendo muito elevado (1,80) o valor desvio-padrão do parâmetro da frequência.

Por sua vez quanto ao modelo c), e em face daquele valor estimado de periodicidade,

começou-se por impôr o valor de 6 e de 7 anos para o ciclo aditivo. Em ambos os casos os valores de R^2 e R_D^2 foram idênticos aos dos modelos anteriores. No entanto, os correlogramas dos resíduos não apresentaram agora o aspecto de ruído branco, já que se destacou um valor significativo a 95% de confiança, correspondente ao desfazamento 3. Procedeu-se, então, à imposição de período igual a 3, tendo-se ultrapassado aquele problema e elevado ligeiramente o R^2 para 0,944 e o R_D^2 para 0,1. No entanto, tal como anteriormente, a ciclicidade encontrada explica valores muitíssimo reduzido da variação do Produto, em escala semelhante à componente irregular.

O PERÍODO 1925-86

Já referimos anteriormente a diferença encontrada antes e depois dos meados dos anos 20. De facto, e embora sublinhemos que os valores encontrados para as 1^{as} décadas são estimados, é de realçar na 1^a fase, p.ex., as seguintes taxas de enorme variação do Produto:

1918 : -25%; 1921 : -20%; 1922 : +26%.

Ignorando o 1^o período, tinha-se concluído a caracterização da série como a de um passeio aleatório + taxa de impulso. A utilização da modelização estrutural veio obviamente confirmar esta perspectiva. Assim obteve-se agora um $R^2 = 0.9978$ e sobretudo um $R_D^2 = 0.0015$, i.e., praticamente nulo.

3 - Conclusões

- O Produto real tem 3 fase distintas quanto à instabilidade do seu crescimento. As mesmas, coincidem aproximadamente com a 1^a República, Estado Novo e Pós - 25 de Abril, não devem levar a sobrevalorizar-se o papel explicativo das mutações do regime, já que há a considerar importantes mudanças também ocorridas na economia internacional.
- Não havendo razão para considerar que o padrão da última fase é o único relevante na extrapolação dos anos vindouros, uma análise a mais prazo (mais de 30 anos) revela que as flutuações cíclicas são responsáveis pela explicação da variância do Produto num montante inferior a cerca de 1%. Esta conclusão é válida para os diferentes modelos ajustados, com excepção da exponencial para todo o período (1913-86), mas parecendo este abusivo, pelo que atrás referimos.
- Quer na tentativa, para os 'ciclos clássicos', de isolar a flutuação cíclica de inflexões da tendência ou da componente irregular, quer mais nitidamente na definição dos 'ciclos de crescimento', a questão primordial da análise estatística parece ser a de como modelizar a tendência-determinística ou estocástica.
- Ainda que a existirem razões da Teoria Económica que perfilhem a 1^a alternativa, a nós parece-nos que a adopção, neste caso, de funções polinomiais no tempo, se ganha em aderência relativamente a funções lineares ou exponenciais, carece de interpretação económica que as imponham face a tendências estocásticas.

Por seu lado; confirma-se que, embora perante valores elevados de aderência (R^2), se podem estar a introduzir nos resíduos, sem aparente significado, grandes fases acima e abaixo da tendência, as quais vêm deturpar a posterior utilização da análise frequencista.

Confirma-se, nesse sentido, as análises de autores como Nelson e Kang (1981), que sublinharam o aparecimento, neste caso, de um pico dominante nas baixas frequências.

- No quadro não determinístico, a análise efectuada com modelos ARIMA e estruturais conclui pela frequência na caracterização do PIB como "Difference-Stationary" e nomeadamente como passeio aleatório + impulso, que se acentua a partir de meados dos anos 20.
- Ainda assim é possível notar-se uma ciclicidade residual, de reduzida amplitude, de período médio igual a 3 anos - flutuação estável que se torna também evidente na taxa de crescimento do PIB. Por sua vez, na série bruta, os valores entre picos ('ciclos clássicos'), apresentam médias de 7-8 anos, mas com uma enorme variabilidade que lhe retira significado.

Referências

- BARBOSA, António Pinto [1985], Inflação e Produção em Portugal 1953-80, "Economia", U.Católica Portuguesa, Vol.I, nº1.
- BEVERIDGE, SE. e NELSON, C.R. (1981), A New Approach to Decomposition of Economic Time Series into Permanent and Transitory Components with Particular Attention to Measurement of the Business Cycle, J.of Monetary Economics, Vol.7, nº2, March, North-Holland and P.C.
- CARTAXO, R.J. e ROSA, N.E.S> [1986], Séries longas para as contas Nacionais Portuguesas 1958-1985, Doc.Trabalho nº15, Banco de Portugal

- HARVEY, A.C. [1984], Trends and Cycles in Macroeconomic Time Series, London School of Economics
- NELSON, C.R. e KANG, H. [1981], Spurious Periodicity in Inappropriately Detrended Time Series, *Econometrica*, Vol.49, nº3, pp.741-751.
- NELSON, C.R. e PLOSSER, C.I. [1982], Trends and Random Walks in Macroeconomic Time Series - Some Evidences and Implications, *Journal of Monetary Economics*, 10, pp.139-162
- PERRON, P. [1987], Trends and Random Walks in Macroeconomic Time Series - Further Evidence from a New Approach, *Journal of Economic Dynamics and Control*, 12 (1988), pp.297-332, NORTH-HOLLAND
- SOARES, João A.Ó. [1989], Métodos Estatísticos de Análise da Ciclicidade em Séries Económicas, IST-U.T.L.
- VALÉRIO, Nuno [1983], O Produto Nacional de Portugal entre 1913 e 1947 - Uma Primeira Aproximação, Separata da Revista de História Económica e Social, Março



DESENVOLVIMENTO DE ALGORITMOS APROXIMATIVOS POR ACERCAMENTO: ESPECIFICAÇÃO FORMAL

Laira Toscani*
Paulo A. S. Veloso**

Resumo: Acercamento é um método de desenvolvimento de algoritmos aproximativos para uma classe de problemas de maximização, que satisfaz um dado requerimento de exatidão. A especificação formal apresentada utiliza tipos abstratos de dados.

Abstract: Rounding is a method for the development of approximate algorithms to a class of maximization problems, that satisfies a given accuracy requirement. The method is formally specified through abstract data types.

1 - Introdução

O Acercamento é um método de desenvolvimento de algoritmos aproximativos (mdaa) para um certo problema de optimização que será definido na secção seguinte. Horowitz ([HOR 78]) chama este mdaa de "Rounding". É um método interessante porque dado um requerimento de exatidão e regras de dominância, gera um algoritmo aproximativo polinomial no tamanho da entrada e inversamente polinomial ao requerimento da exatidão. Isto é, o algoritmo resultante é polinomial e gera soluções "satisfatoriamente" próximas da solução óptima.

A especificação formal do método consiste de um programa abstrato e um conjunto de axiomas. O programa abstrato mostra a estratégia do método. Os axiomas definem o inter-relacionamento das funções e dos predicados e permitem a verificação da correção do programa abstrato. A abstração é desenvolvida em dois níveis.

Para aplicar esta estratégia a uma instância do problema é suficiente escrever um módulo de implementação definindo as funções abstratas em termos do problema. Com isto, ter-se-á um algoritmo aproximativo, que é uma instância do programa abstrato, que satisfaz o requerimento de exatidão $\frac{n}{\xi}$, em que n é o tamanho da instância.

*Doutora em Informática (PUC/RJ, 88); área de interesse: desenvolvimento e complexidade de algoritmos; Profª adjunta da UFRGS-Brasil; atualmente na Universidade Nova de Lisboa, Fac. de Ciências e Tecnologia, Grupo de Programação em Lógica e Inteligência Artificial, Quinta da Torre, 2825 Monte da Caparica, Portugal.

**Ph.D. em Ciência da Computação (Univ. da Califórnia/Berkeley, 75); área de interesse: teoria e metodologia de programação; Prof. Associado na PUC/RJ; Rua Marquês de São Vicente 225, 22.453, Rio de Janeiro, Brasil.

2 - Definição do Problema

Problema: calcular $\max_x \sum_{i=1}^n p_i x_i$ (função objectivo) restrito a $\sum_{i=1}^n a_{ij} x_i \leq b_j$, $a_{ij} \geq 0$ e $1 \leq j \leq n$; $x_i = 0$ ou 1 , $p_i \geq 0$ e $1 \leq i \leq n$.

Definições:

Seja $I=(P, A, B)$, onde $A=(a_{ij})$ é uma matriz $n \times n$, $B=(b_1, b_2, \dots, b_n)$, $P=(p_1, p_2, \dots, p_n)$ e $a_{ij}, b_j, p_j \geq 0$, $1 \leq i, j \leq n$. Seja $s = (s_1, s_2, \dots, s_i)$, onde $s_j = 0$ ou 1 , $1 \leq j \leq i \leq n$.

- I é uma instância do problema.

- n é o tamanho de I .

- s é uma solução parcial possível de I . Se $i=n$, então s é uma solução possível e se além

disso $\sum_{j=1}^i a_{jk} s_j \leq b_k$, $1 \leq k \leq i$, então s é uma solução parcial viável de I .

- $FO(I, s, i) = \sum_{j=1}^i p_j s_j$. Se s é a solução óptima de I restrito a i variáveis, $FP^*(I, i) =$

$= \sum_{j=1}^i p_j s_j$ e $FP^*(I, n) = F^*(n) = \sum_{j=1}^n p_j s_j$, em que s é a solução óptima para I .

- $s^{(k)}$ com $1 \leq k \leq n$ é o conjunto das soluções viáveis de I restrito às k primeiras variáveis

3 - Definição do Método de Acercamento

O método consiste em transformar a instância original I em uma instância I' , cuja solução está próxima da solução óptima da instância original, transformar I' em uma instância I'' , que acelera o processo de construção de $S^{(0)}, S^{(1)}, \dots, S^{(n)}$, sem alterar a solução óptima. Em $S^{(n)}$ é calculada a solução óptima de I'' (e de I') e transportada para I .

Características do Método

Seja I uma instância dada de tamanho n , I' e I'' as instâncias calculadas pelo método, ξ o requerimento de exatidão, s a solução de I'' encontrada através do método e s^* a solução óptima de I . Então:

- a solução encontrada está satisfatoriamente próxima da solução óptima, isto é, $|F^*(I) - F^*(I'')| / F^*(I) \leq \xi$;

- a solução encontrada s é a solução óptima das instâncias I' e I'' ;

- os métodos de transformação de I em I' e de I' em I'' tem complexidade linear no comprimento de I .

- a cardinalidade do conjunto de soluções parciais $S^{(k)}$ é da ordem $k(n/\xi)$, com $k = 1, 2, \dots, n$. O que vai garantir uma complexidade total do método polinomial em n/ξ .

4 - Especificação Formal do Método

A especificação formal é constituída do programa abstrato e do conjunto de axiomas.

Definição dos Domínios:

PO - domínio de instâncias do problema. Se $I \in PO$ então $I = (P, A, B)$, onde $A = (a_{ij})$ é uma matriz $n \times n$, $B = (b_1, b_2, \dots, b_n)$ e $P = (p_1, p_2, \dots, p_n)$.

\mathcal{R} - domínio das entradas adicionais x e Lb ; x é o requerimento de exatidão; Lb estimativa inferior para o valor óptimo da função objectivo.

RD - domínio das regras de dominância de soluções viáveis.

SP - domínio das soluções parciais: conjunto de conjuntos de i -uplas ($0 \leq i \leq n$), soluções possíveis do problema restrito às i primeiras variáveis.

SL - domínio das soluções: conjunto das n -uplas, soluções possíveis para o problema.

\mathcal{R} - contradomínio da função tamanho, que define o tamanho da instância do problema.

V - domínio dos valores possíveis da função objectivo.

Funções Auxiliares:

- tamanho: $PO \rightarrow \mathcal{R}$, calcula o tamanho da instância do problema considerado.

- transformal: $PO \times \mathcal{R}^2 \rightarrow PO$, transforma uma instância I , considerando Lb e ξ , em

- uma nova instância I' cuja solução está próxima da solução da instância original I .
- transforma2: $PO \times \mathcal{R}^2 \rightarrow PO$, transforma a instância I' em uma nova instância I'' com mesma solução.
- inicializa: $X^0 \rightarrow SP$, inicializa o conjunto de soluções viáveis para o caso de zero variáveis ($X^0 = \{\epsilon\}$, ϵ string vazio).
- calcula: $PO \times SP \times RD \times \mathcal{R} \rightarrow SP$, calcula o conjunto de soluções parciais viáveis, obtido considerando o conjunto de soluções parciais dadas, incluindo uma variável a mais no problema considerado e fazendo as eliminações permitidas pelas regras de dominância.
- calculaSO: $PO \times SP \rightarrow SL$, escolhe entre os valores oferecidos no conjunto de soluções viáveis a melhor solução para o problema.
- calculaR: $PO \times L \rightarrow V$, calcula a função objectivo para a solução considerada (calculaR(I, S) = FO(I, s , tamanho(I))).
- F*: $PO \rightarrow V$ t.q. $F^*(I) = FO(I, s, tamanho(I))$, onde s = solução óptima de I .
- cardinalidade: $SP \rightarrow \mathcal{R}$, t.q. cardinalidade(S) = cardinalidade do conjunto S .

Predicados Auxiliares:

- Bemtransf: $PO^2 \times \mathcal{R} \rightarrow \{T, F\}$, Bemtransf(I, I', ξ) = $|F^*(I) - F^*(I')| / F^*(I) \leq \xi$.
- MesmaSol: $PO^2 \rightarrow \{T, F\}$, MesmaSol(I, I') = $(\forall s) [SolOtima(I, s) \Leftrightarrow SolOtima(I', s)]$.
- SolOtima: $PO \times SL \rightarrow \{T, F\}$, SolOtima(I, s) = s é a solução óptima de I .
- SolAprox: $PO \times SL \times \mathcal{R} \rightarrow \{T, F\}$, SolAprox(I, s, ξ) = $|F^*(I) - FO(I, s, tamanho(I))| / F^*(I) \leq \xi$.
- LimiteInf: $PO \times \mathcal{R} \times \{T, F\}$, LimiteInf(I, Lb) = $F^*(I) \leq Lb$.
- ContemSol: $PO \times SP \times \mathcal{R} \rightarrow \{T, F\}$, contemSol(I, S, k) = S contém uma solução óptima de I , restrita a k variáveis.
- BemDef: $PO \times RD \rightarrow \{T, F\}$, BemDef(I, RD) = $(\forall s_1, s_2) \{ [s_1, s_2 \in SL \wedge Domina(s_1, s_2, RD) \wedge SolOtima(I, s_2)] \Rightarrow SolOtima(I, s_1) \}$.
- CardBoa: $SP \times \mathcal{R}^2 \times \mathcal{R} \rightarrow \{T, F\}$, CardBoa(S, k, n, ξ) = $(\exists C_1, C_2 \in \mathcal{R})$ (cardinalidade(S) $\leq C_1 + C_2 \cdot k \cdot (n/\xi)$).
- Domina: $SL^2 \times RD \rightarrow \{T, F\}$, Domina(s_1, s_2, Rd) = s_2 domina s_1 pela regra de dominância Rd .

Axiomas: ACAX

- AC1: $(\forall I) (\forall \xi) (\forall Lb) \{ LimiteInf(I, Lb) \Rightarrow Bemtransf(I, transformal(I, \xi, Lb)) \}$
- AC2: $(\forall I) (\forall \xi) (\forall Lb) \{ MesmaSol(I, transformal(I, \xi, Lb)) \}$
- AC3: $(\forall I) \{ ContemSol(I, inicializa(I, 0)) \}$
- AC4: $(\forall I) (\forall S) (\forall k) (\forall Rd) \{ BemDef(I, Rd) \Rightarrow [ContemSol(I, S, k) \Rightarrow ContemSol(I, calculaS(I, S, Rd, k), k+1)] \}$
- AC5: $(\forall I) (\forall S) (\forall \xi) \{ ContemSol(I, S, tamanho(I)) \Rightarrow SolOtima(I, calculaSO(I, S)) \}$
- AC6: $(\forall I) (\forall I') (\forall I'') (\forall \xi) (\forall s) \{ [Bemtransf(I, I', \xi) \wedge MesmaSol(I', I'') \wedge SolOtima(I'', s)] \Rightarrow SolAprox(I, s, \xi) \}$
- AC7: $(\forall I) (\forall \xi) \{ CardBoa(inicializa(I, 0), tamanho(I), \xi) \}$
- AC8: $(\forall I) (\forall S) (\forall Rd) (\forall k) \{ BemDef(I, Rd) \Rightarrow [CardBoa(S, k, tamanho(I), \xi) \Rightarrow CardBoa(calculaS(I, S, Rd), k+1, tamanho(I), \xi)] \}$

Programa: ACERCAMENTO

entrada: I, ξ, Lb, Rd

1. $I' \leftarrow transformal(I, \xi, Lb)$;
2. $I'' \leftarrow transformal(I', \xi, Lb)$;
3. $S \leftarrow inicializa(I)$;
4. $n \leftarrow tamanho(I)$

5. para $k = 0$ até $n - 1$ faça
 6. $S \leftarrow \text{calcula}S(I'', S, Rd, k)$;
 7. fim-para
 8. $s \leftarrow \text{calcula}SO(I'', S)$;
 9. $v \leftarrow \text{calcula}R(I, s)$;
- saída: s, v

Entradas:

- $I \in PO$: instância de entrada.
- $\xi \in \mathcal{R}$: requerimento de exatidão.
- $Rd \in RD$: regras de dominância.

Saídas:

- $s \in SL$: solução alcançada.
- $v \in V$: valor da função objectivo para a solução alcançada.

A prova da correção do programa ACERCAMENTO pode ser extensa e cansativa, mas não é difícil e pode ser encontrada no Apêndice E1 de [TOS 88].

5 - Uma Implementação

Uma implementação de um tipo abstrato é uma definição do tipo em um nível de abstração mais baixo. Uma implementação do tipo definido na secção anterior pode ser uma definição algorítmica das funções *transformal*, *transforma2*, e *calculas*.

Definição de transformal

Programa: *transformal*

entrada: $(P = (p_1, \dots, p_n), A, B, \xi, Lb)$

1. $n \leftarrow \text{tamanho}(P)$
 2. para $i = 1$ até n faça
 3. $q_i \leftarrow \lfloor \frac{p_i \cdot n}{Lb \cdot \xi} \rfloor \frac{Lb \cdot \xi}{n}$
 4. fim-para
 5. $Q \leftarrow (q_1, q_2, \dots, q_n)$
- saída: (Q, A, B)

Definição de transforma2

Programa: *transformal*

entrada: $(p = (p_1, \dots, p_n), A, B), \xi, Lb$

1. $n \leftarrow \text{tamanho}(p)$
 2. para $i=1$ até n faça
 3. $q_i \leftarrow \frac{p_i \cdot n}{Lb \cdot \xi}$
 4. fim-para
 5. $Q \leftarrow (q_1, q_2, \dots, q_n)$
- saída: (Q, A, B)

Definição de calculas

$SS = \bigcup_{S \in SP} =$ conjunto de todas i -uplas $1 \leq i \leq n$, soluções parciais

Funções e Predicados Auxiliares

- *completa*: $SS \times \mathcal{K} \leftarrow SS$, *completa* $((x_1, x_2, \dots, x_k), k)$ considera a k -ésima variável do problema resultando em uma solução possível com mais um elemento, $(x_1, x_2, \dots, x_k, 1)$.
- *Solviável* $SS \times SP \times PO \times \mathcal{K} \leftarrow \{T, F\}$, *Solviável* (s, S, I, k) verifica se a solução parcial s é viável para k variáveis.
- *restringe*: $SP \times RD \times PO \rightarrow SP$, *restringe* (S, Rd, I) restringe o conjunto de soluções parciais S de I , de maneira que o novo conjunto de soluções parciais não contenha duas soluções s_1, s_2 , com s_1 dominando s_2 (dominando de acordo com as regras de dominância Rd).
- *Domina*: $SS^2 \times RD \rightarrow \{T, F\}$, *Domina* (s_1, s_2, Rd) verifica se s_1 domina s_2 de acordo com as regras de dominância Rd .

- Éviável: $PO \times SL \rightarrow \{T, F\}$, Éviável $(I, s) = \sum_{i=1}^n a_{ij}s_i \leq b_j$, onde $n = \text{tamanho}(I)$, $I = (P, A, B)$, $A = (a_{ij})$, $B = (b_1, b_2, \dots, b_n)$ e $s = (s_1, \dots, s_n)$.

Programa: calculas

$\{\varphi(I, S, Rd, k) = \text{BemDef}(I, Rd) \wedge \text{ContemSol}(I, S, K)\}$

entrada: I, S, Rd, k

1. $S' \leftarrow \phi$
2. para $s \in S$ faça
3. $s' \leftarrow \text{completa}(s, k)$
4. se Solviável $(s', S, I, k+1)$ então $S' \leftarrow U \{s'\}$
5. fim - para
(invariante: $\text{ContemSol}(I, S', k+1) \wedge \text{BemDef}(I, Rd)$)
6. $S \leftarrow \text{restringe}(S', Rd, I)$
saída: S
 $\{\psi(I, S, Rd, k) = \text{ContemSol}(I, S, k+1)\}$

Axiomas: CalSAX

- ACS1. $(\forall I) (\forall S) (\forall S') (\forall k) \{[\text{ContemSol}(I, s, k) \wedge S' = U \{\text{completa}(s, k)\} \wedge \text{Solviável}(\text{completa}(s, I), S, I)] \Rightarrow \text{ContemSol}(I, S', k+1)\}$
- ACS2. $(\forall I) (\forall S) (\forall k) (\forall rd) \{[\text{Bemdef}(I, Rd) \wedge \text{ContemSol}(I, S, k)] \Rightarrow [\text{ContemSol}(I, \text{restringe}(S, Rd, I), k) \wedge (\forall s_1, s_2 \in \text{restringe}(S, Rd, I)) (\neg \text{Domina}(s_1, s_2, Rd))]\}$
- ACS3. $(\forall I) (\forall S) (\forall S') (\forall k) \{[\text{ContemSol}(I, S, k)] [(\forall s \in S) \text{Solviável}(\text{completa}(s, k), S, I, k+1) \Rightarrow \text{completa}(s, k) \in S'] \Rightarrow \text{ContemSol}(I, S', k+1)\}$

Note que calculaS, diferente de transforma1 e transforma2, contém funções abstratas cujas semânticas precisam ser definidas e são definidas pelos axiomas CalSAX. Um novo módulo de implementação pode definir as funções abstratas de calculaS (completa e restringe).

A correção da implementação é facilmente provada ([TOS 88], secção 6.2.2).

6 - Exemplo

Este exemplo foi retirado de [HOR 78] e é a solução de uma instância do problema da mochila.

O problema da mochila é uma variante do problema de maximização estudado e pode ser definido assim: c alcular

$$\max_x \sum_{i=1}^n p_i x_i \text{ restrito a } \sum_{i=1}^n w_i x_i \leq M, x_i = 0 \text{ ou } 1, 1 \geq i \leq n.$$

Chame $r = \sum_{j=1}^i p_j x_j$ e $t = \sum_{j=1}^i w_j x_j$, (r, t) são as atribuições possíveis. A regra de

dominância é: (r_1, t_1) domina (r_2, t_2) sss $t_1 \leq t_2$ e $r_1 \geq r_2$.

A instância é $I = (P, W, M)$, onde $P = (p_1, p_2, p_3, p_4, p_5) = W = (w_1, w_2, w_3, w_4, w_5) = (1, 2, 10, 100, 1000)$, $M = 1112$. Para esta instância $r = t$.

Em $V^{(i)}$ serão representados os valores da função objectivo para as soluções parciais viáveis, ao invés das próprias soluções parciais viáveis.

$$V^{(0)} = \{0\},$$

$$V^{(1)} = \{0, 1\},$$

$$V^{(2)} = \{0, 1, 2, 3\},$$

$$V^{(3)} = \{0, 1, 2, 3, 10, 11, 12, 13\},$$

$$V^{(4)} = \{0, 1, 2, 3, 10, 11, 12, 13, 100, 101, 102, 103, 110, 111, 112, 113\},$$

$$V^{(5)} = \{0, 1, 2, 3, 10, 11, 12, 13, 100, 101, 102, 103, 110, 111, 112, 113, 1000, 1001, 1002, 1003, 1010, 1011, 1012, 1013, 1100, 1101, 1102, 1103, 1110, 1111, 1112\}$$

A solução óptima é (0,1,1,1) com valor óptimo 1112. Usando o método de Acercamento, com $\xi = (1/10)$, $Lb = \max \{p_i\} = 1000$, tem-se $I'' = (Q, W, M)$, $Q = (0, 0, 0, 5, 50)$,

$$q_i = \lfloor \frac{P_i \cdot n}{Lb \cdot \xi} \rfloor = \lfloor \frac{P_i}{20} \rfloor.$$

$$V(0) = V(1) = V(2) = V(3) = \{(0,0)\}$$

(Obs.: (0,1) não está em $V(3)$ porque é dominado por (0,0)).

$$V(4) = \{(0, 0), (5, 100)\}, V(5) = \{(0, 0), 5, 100\}, (50, 1000), (55, 1100)\}.$$

A solução óptima em I'' é $s = (0, 0, 0, 1, 1)$, com valor $FO(I'', s, 5) = 55$. s transportada para I resulta $FO(I, s, 5) = 1100$. $F^*(I) - FO(I, s, 5) / F^*(I) = (1112 - 1100) / 1112 = 12 / 1112 < 0.011 < \xi$.

7 - Conclusão

Muitas vezes os dados de um problema são aproximados, assim uma solução aproximada "suficientemente" da solução exata é tão significativa quanto a própria solução exacta. Baseados nesse facto foram desenvolvidos métodos de solução de problemas que encontram soluções aproximadas, isto é, algoritmos aproximativos. A condição "suficientemente próxima", para muitas aplicações é uma exigência de qualidade mínima para aceitação de uma solução aproximada. Esta exigência é posta neste texto como requerimento de exactidão.

Métodos de desenvolvimento de algoritmos como Divisão e Conquista e Programação Dinâmica já foram exaustivamente estudados [VEL 80], [TOS 85], [TOS 86], [TOS 88]. Métodos para desenvolvimento de algoritmos aproximativos são mais complexos, já que os algoritmos gerados tem que satisfazer um certo requerimento de exactidão e de complexidade, pois um algoritmo aproximativo para ser útil tem que dar uma solução satisfatoriamente próxima da solução exacta em tempo significativamente menor que o algoritmo exacto.

A especificação formal do método de Acercamento pretende ser uma contribuição no sentido de aumentar o conhecimento sobre o método, melhorando a compreensão do mesmo e facilitando a sua utilização. Além disso, a especificação formal necessária para um estudo teórico do método, como por exemplo a análise comparativa entre mdaa's.

Referências

- [1] Horowitz, E. & SAHNI, S. - *Fundamentals of Computer Algorithms*, Potomac, Md. Comp. Sci. Press, 1978.
- [2] Toscani, L.V. & Veloso, P.A.S. - *Uma especificação Formal para a Programação Dinâmica*. In: Seminário Integrado de Software e Hardware, 12., Porto Alegre, 20-27/Jul, 1985. Anais. Porto Alegre, SBC/CLEI/UFRGS, 1985. p. 477-86.
- [3] Toscani, L.V. & Veloso, P.A.S. - *Divisão e Conquista: análise da complexidade*. In: Seminário Integrado de Software e Hardware, 13, Olinda, 19-25/Jul., 1986. Anais. Recife SBC, 1986. p. 89-104. Journal SIAM 12 (1964) 663-665.
- [4] Toscani, L.V. - *Métodos de Desenvolvimento de Algoritmos : Análise Comparativa e de Complexidade*, Tese de Doutorado, Rio de Janeiro, Depto. de Informática, Purc/RJ. 1988. (1977) 173-194.
- [5] Veloso, P.A.S. - *Divide-and-Conquer via data types* In: Latin American Conference on Informatics, 7, Caracas, 1980. Proceedings.

FUZZY LINEAR PROGRAMMING - A Tentative Survey

by

A. C. Rosa
Departamento de Matemática
Universidade de Coimbra

J. C. N. Clímaco
Departamento de Engenharia Electrotécnica
Universidade de Coimbra

Abstract: One of the areas in which fuzzy sets have been applied most extensively is in the modelling of managerial decision making. This paper is intended to present the state of the art concerning the application of fuzzy sets to linear programming. Two classes of approaches are considered: the flexible fuzzy linear programming ones and the robust fuzzy linear programming approaches. Its bibliographic support ends in May 1986.

General Introduction: Probability theory is often used to deal with the imprecision in the formulation and solution of systems and decision problems. However, it is becoming increasingly clear, that in the case of many real world problems involving large scale systems, the major source of imprecision should more properly be labeled as "fuzziness" rather than "randomness". By fuzziness, we mean the type of imprecision which is associated with the lack of sharp transition from membership to nonmembership, as in "tall men", "small numbers", etc.

Since its inception 20 years ago the theory of fuzzy sets has advanced in a variety of ways and in many disciplines. Theoretical developments and applications of this theory can be found in many domains. Publications are already very specialized and widely scattered over many areas.

In general, whenever linear programming models are used for dealing with real world problems several sources of imprecision are neglected. It is impossible describing exactly the physical real world problems using traditional mathematical models, because of unexpected relevant events, data value variations, linearization of truly nonlinear phenomena, subjective human judgements, etc. On the other hand, the limited precision of computational calculation techniques and the analysis and evaluation of the solutions (since the model fits the reality only partially) are also sources of imprecision. During the last few years the researchers have tried to include, as much as possible, these aspects in linear programming models. Three types of approaches must be referred: interval linear programming, stochastic linear programming and fuzzy linear programming. In stochastic linear programming it is implicitly supposed that the imprecision is due to random phenomena, and it uses the probability theory axiomatic system. It must be emphasized that the roots of the subjective imprecision are completely different. A dichotomic modelling language is not adequate when the nondeterministic factors are related to subjective human judgements. Roughly speaking, the language is the interface among the human thinking, the object system and the model itself. This means that when one intends the incorporation of this kind of imprecision in a decision model a semantic language is adequate. The fuzzy decision models enable the consideration of fuzzy variables like "reasonable profits", "promising merchandizing", "insufficient production", etc, intrinsically related to the natural language features.

In this paper we deal with a special class of fuzzy decision models, the fuzzy linear programming approaches.

We have tried to describe the subject as detailed as necessary so that it may be comprehended by those who have not been exposed to this theory. Examples are used for illustrating the concepts more clearly. Numerous references are included for the interested reader.

In part 1, basic definitions of fuzzy sets and algebraic operations are presented which will support further considerations. Some extensions are then presented by introducing additional concepts and alternative operators. The extension principle is stated and fuzzy arithmetic is introduced. The coverage of fuzzy concepts in part 1 is restricted to those which are directly used in the remaining parts of the text.

In part 2, the conceptual framework for optimization is reviewed and then particularized for fuzzy linear programming. The concept of fuzzy decision, as defined by Bellman and Zadeh [1] and the 'fuzzification' of a linear program are presented: Flexible and Robust programming are also introduced.

Part 3 is dedicated to Flexible Programming. It assumes that the constraints and/or the objective are not precisely formulated and states that the vagueness results from formulations of the type "approximately equal", "as close as", "much bigger", etc, while keeping non-fuzzy the parameters of the model.

Finally Part 4 describes Robust Programming by studying linear programs where coefficients are fuzzy numbers. Two points of view are considered based on the comparison between fuzzy numbers and on the usual 'inclusion' between fuzzy sets, respectively.

1. - Introduction to the Fuzzy Sets Theory

1.1. - Basic Definitions

Let X be a collection of objects generically designated by x .

In the classic theory of sets, the belonging relationship of one element related to a subset A of X is frequently defined by the corresponding characteristic function ($1_A(x) = 1$ iff $x \in A$ and $1_A(x) = 0$ iff $x \notin A$).

The fuzzy set's notion, innovated by Zadeh [35], abandons the binary character of this relation, enabling that an object is only partially a member of a set.

Def 1.1 - A fuzzy set \tilde{A} , of X , is a set of ordered pairs, $\tilde{A} = \{(x, \mu_{\tilde{A}}(x)), x \in X\}$, with $\mu_{\tilde{A}} : X \rightarrow [0,1]$.

Obs 1.1 :

- (i) $\mu_{\tilde{A}}$ is the membership function of the fuzzy set \tilde{A} and $\mu_{\tilde{A}}(x)$ expresses the membership's level (or compatibility's level) of x with respect to \tilde{A} ;
- (ii) all classic subset of X will be denominated as rigid, the membership and characteristic function, coinciding, in this case;
- (iii) the empty set is defined by $\mu_{\emptyset}(x) = 0, \forall x \in X$ and we obviously obtain $\mu_X(x) = 1, \forall x \in X$;
- (iv) several alternative notations for the fuzzy set representation can be found in the literature [9,23,36,41].
- (v) we designate by $\tilde{\mathcal{P}}(X)$ the set of fuzzy subsets in X .

Let $\tilde{A}, \tilde{B} \in \tilde{\mathcal{P}}(X)$.

Def 1.2 - The support of \tilde{A} , $\mathcal{S}(\tilde{A})$, is the rigid subset of $X : \{x \in X : \mu_{\tilde{A}}(x) > 0\}$.

Def 1.3 - We designate the height of \tilde{A} as the lower upper bound of $\mu_{\tilde{A}}(x)$, i.e.,

$$alt(\tilde{A}) = \sup_{x \in X} \mu_{\tilde{A}}(x)$$

Def 1.4 - The set \tilde{A} is normalized if there exists $x \in X$ such that $\mu_{\tilde{A}}(x) = 1$.

Ex 1.1 - Let $X = \{50, 60, 70, 80, 90, 100, 110, 120, 130, 140\}$ be the medium speed set (km/h) for the cars that travel long distances. Let us consider the fuzzy set \tilde{A} of the "comfortable medium speeds for long distances".

\tilde{A} can then be defined (for a determined individual) by:

$$\tilde{A} = \{(50, 0.2), (60, 0.5), (70, 0.8), (80, 0.9), (90, 1.0), (100, 0.7), (110, 0.6), (120, 0.3)\}.$$

Def 1.5 - We say that $\bar{A} = \bar{B}$ iff $\mu_{\bar{A}}(x) = \mu_{\bar{B}}(x), \forall x \in X$

Def 1.6 - We say that : $\bar{A} \subseteq \bar{B}$ iff $\mu_{\bar{A}}(x) \leq \mu_{\bar{B}}(x), \forall x \in X$;

$\bar{A} \subset \bar{B}$ iff $\mu_{\bar{A}}(x) < \mu_{\bar{B}}(x), \forall x \in X$

(inclusion in a wide sense and a narrow sense, respectively).

Def 1.7 - Let $\alpha \in \mathbb{R}$. The rigid set $A_\alpha = \{x \in X : \mu_{\bar{A}}(x) \geq \alpha\}$ is designated as α level cut or α level set of \bar{A} .

Obs 1.2 :

- (i) If the inequality in the wide sense is substituted for one strict inequality, the set takes the designation of a strong α level cut of \bar{A} ;
- (ii) We have the following properties of the α level cut:
 - $A_\alpha = X, \forall \alpha \leq 0$ and $A_\alpha = \emptyset, \forall \alpha > 1$;
 - The succession of level cuts of the set \bar{A} is non-increasing rising (in the sense of inclusion);
- (iii) The membership function of a fuzzy set can be expressed in terms of the characteristic functions of it's α level cuts [9,36]:

$$\mu_{\bar{A}}(x) = \sup_{\alpha \in (0,1)} \min(\alpha, \mathbf{1}_{A_\alpha}(x)), \forall x \in X.$$

1.2. - Operations on fuzzy sets

Next we present the generalization of the classic union, intersection and complementarity operations proposed by Zadeh [35].

In what follows, we will suppose \bar{A} and \bar{B} are fuzzy sets in X (fig. 1.1).

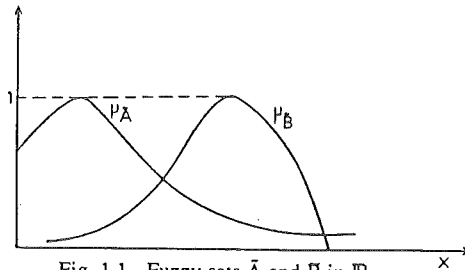


Fig. 1.1 - Fuzzy sets \bar{A} and \bar{B} in \mathbb{R}

Def 1.8 - The intersection's membership function, $C = \bar{A} \cap \bar{B}$, is defined by:

$$\mu_C(x) = \min(\mu_{\bar{A}}(x), \mu_{\bar{B}}(x)), \forall x \in X \text{ (fig.1.2)}.$$

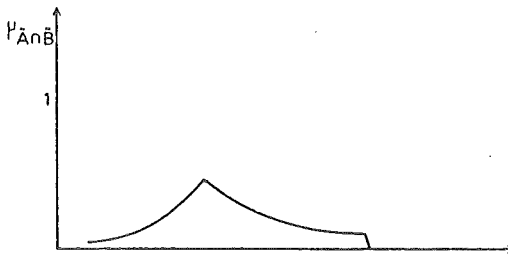


Fig. 1.2 - Membership function of $\bar{A} \cap \bar{B}$

Def 1.9 - The reunion's membership function, $D = \bar{A} \cup B$, is pointedly defined by:

$$\mu_D(x) = \min(\mu_{\bar{A}}(x), \mu_B(x)), \forall x \in X \text{ (fig.1.3).}$$

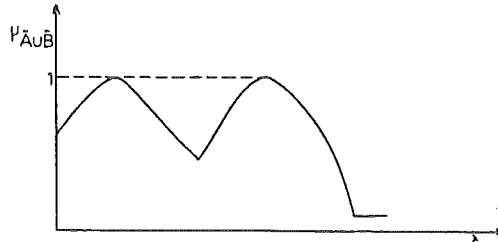


Fig. 1.3 - Membership function of $\bar{A} \cup B$

Def 1.10 - We designate the complement of A as \bar{A}^c . This is the fuzzy set in X, whose membership function is given by:

$$\mu_{\bar{A}^c}(x) = 1 - \mu_A(x), \forall x \in X \text{ (fig.1.4).}$$

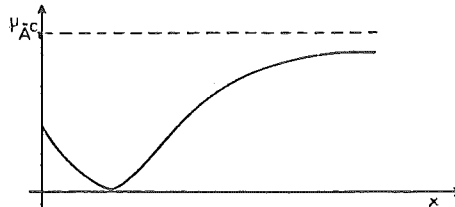


Fig. 1.4 - Membership function of \bar{A}^c

Obs 1.3 :

- (i) Note that, when \bar{A} and B are not fuzzy ($\mu_{\bar{A}} \equiv 1_A$ and $\mu_B \equiv 1_B$), the expressions presented allow the obtention of the usual operations of the classic theory of sets;
- (ii) The defined operations verify all properties of the classic theory of sets [9], with the exception of the half excluded's law, once in general, $\bar{A} \cap \bar{A}^c \neq \emptyset$ and $\bar{A} \cap \bar{A}^c \neq X$;
- (iii) The justification of the choice of the max (min) operators for characterizing the

reunion (intersection) operation in $\tilde{\mathcal{P}}(X)$ was made, in axiomatic form, by *Bellman* and *Giertz* [1]. Their argument is based on a logical point of view, interpreting the intersection as "logical and", the reunion as "logical or" and the fuzzy set \bar{A} as the "statement" the element x belongs to the set \bar{A} ". The same authors showed that is more difficult to justify the def. 1.10 in order to translate the complementar notion of a fuzzy set.

1.3 - Complementary topics

Without the intention of being exhaustive, we will indicate, on one hand, some extensions of the introduced concepts and on the other hand, we will make reference to other concepts that seem important to us in order to develop a study on this theory.

1.3.1. - Extensions

With respect to the definition, we can admit that the attribution values of the considered objects are fuzzy sets in [0,1] (fuzzy set type 2 notion [9,41]), or also that the value of the membership function, for $x \in X$, is a random variable (probabilistic set notion [9,41]). Finally, a new extension can be obtained by considering that the arrival set of the membership function of a fuzzy set is more general than the interval [0,1] (L-fuzzy set's notion [9]).

The first two arise in an attempt to include uncertainty in the primitive definition mainly due to the lack of precision of the attribution values in real situations [9]. For this reason, other relations of inclusion and equality between fuzzy sets were also introduced.

Nevertheless, it is on the domain of the reunion and intersection operations established in $\tilde{\mathcal{P}}(X)$ that the more diverse suggestions can be found [9]. Such suggestions differ not only by the shape and degree by which the respective justifications are elaborated (intuitive argumentation, empirical or axiomatic justification), but also by the character of generality and adaptation (parametric families of operators, classes satisfying certain properties).

The mathematic representation of the intersection operation, may be:

- (i) the algebraic product, $\tilde{A} \cdot \tilde{B}$ ($\mu_{\tilde{A} \cdot \tilde{B}} = \mu_{\tilde{A}} \cdot \mu_{\tilde{B}}$) [36];
- (ii) the limited difference, $\tilde{A} \ominus \tilde{B}$ ($\mu_{\tilde{A} \ominus \tilde{B}} = \max(0, \mu_{\tilde{A}} + \mu_{\tilde{B}})$)[41];
- (iii) the Hamacher's intersection operator, that define $\tilde{A} \subset \tilde{B}$ in the following manner:

$$\mu_{\tilde{A} \cup \tilde{B}} = \frac{\mu_{\tilde{A}} \cdot \mu_{\tilde{B}}}{y + (1-y)(\mu_{\tilde{A}} + \mu_{\tilde{B}} - \mu_{\tilde{A}} \cdot \mu_{\tilde{B}})}, y \geq 0.$$

In the same way, it was suggested for the fuzzy sets reunion:

- (i) the probabilistic sum, $\tilde{A} + \tilde{B}$ ($\mu_{\tilde{A} + \tilde{B}} = \mu_{\tilde{A}} + \mu_{\tilde{B}} - \mu_{\tilde{A}} \cdot \mu_{\tilde{B}}$);
- (ii) the limited sum, $\tilde{A} \oplus \tilde{B}$ ($\mu_{\tilde{A} \oplus \tilde{B}} = \min(1, \mu_{\tilde{A}} + \mu_{\tilde{B}})$) [41];
- (iii) the Hamacher's reunion operator, by which:

$$\mu_{\tilde{A} \cap \tilde{B}} = \frac{(1-y)\mu_{\tilde{A}} \cdot \mu_{\tilde{B}} + y(\mu_{\tilde{A}} + \mu_{\tilde{B}})}{y + \mu_{\tilde{A}} \cdot \mu_{\tilde{B}}}, y \geq 0.$$

Hamacher has presented an axiomatic justification, referred by several authors [40,41].

Such diversity of operators for defining the referred operations seems reasonable in comparison with the unity found in the classical theory of sets, since many of those operators (for example the min operator and the algebraic product) lead to the same results if the attributed values are restricted to the values 0 and 1, the same not happening for a wider range of variation.

1.3.2 - Additional notions

One of the basic questions of this theory is to know in what way it is possible to extend, the classic mathematic structures.

The answer can be found in the principle of the extension, proposed by Zadeh [36], that supplies a global and systematic method of generalization of the non-fuzzy notions, leading, namely, to a concept of fuzzy function [41].

1.3.2.1. The principle of extension

Def. 1.11 - Let $\tilde{A}_1, \dots, \tilde{A}_n$ be fuzzy sets in X_1, \dots, X_n respectively. The cartesian product of $\tilde{A}_1, \dots, \tilde{A}_n$, $\tilde{A}_1 \times \dots \times \tilde{A}_n$, is a fuzzy set in the product universe $X_1 \times \dots \times X_n$, whose membership function is defined by :

$$\mu_{\tilde{A}_1 \times \dots \times \tilde{A}_n}(x) = \min_{i=1, \dots, n} \mu_{\tilde{A}_i}(x_i)$$

for all $x = (x_1, \dots, x_n)$ in $X_1 \times \dots \times X_n$.

Def. 1.12 - Let \tilde{A}_i be a fuzzy set in X_i ($i=1, \dots, n$) and f an application of $X = X_1 \times \dots \times X_n$ on the universe Y . Then, starting from the fuzzy set $\tilde{A}_1 \times \dots \times \tilde{A}_n$ we say that f induces to a fuzzy set $\tilde{B} = \{(y), \mu_{\tilde{B}}(y), y \in Y\}$ in Y whose membership function is defined by:

$$\mu_{\tilde{B}}(y) = \begin{cases} \sup_{(x_1, \dots, x_n) \in f^{-1}(y)} \mu_{\tilde{A}_1 \times \dots \times \tilde{A}_n}(x_1, \dots, x_n) & \text{if } f^{-1}(y) \neq \emptyset \\ 0 & \text{otherwise} \end{cases}$$

where $f^{-1}(y) = \{(x_1, \dots, x_n) : y = f(x_1, \dots, x_n)\}$.

Obs. 1.4

We find in [9] some alterations to the principle, namely:

- the substitution of sup by the probablistic sum, combining the probablistic and fuzzy approximation and assuming the dependence of $\mu_{\tilde{g}}(y)$ on the number of (x_1, \dots, x_n) such that $y = f(x_1, \dots, x_n)$.
 - the substitution of the min operator by the product operator, assuming implicitly some interactivity or compensation between A_1, \dots, A_n .
- Another important application of this principle is in the domain of the real fuzzy algebra.

1.3.2.2. Introduction to fuzzy real algebra

We will consider in this section a very special class of fuzzy sets which generalise the classic notion of real number.

The concept of fuzzy real number is due to *Dubois* and *Prade* [6, 7, 9]. In relation to this entity, the authors defined the basic operations (addition, subtraction, ...) and analysed the algebraic structure of the fuzzy real numbers set when provided with such operations.

Def. 1.13 - A fuzzy real number, \tilde{n} , is a fuzzy subset of the straight real line, whose membership function $\mu_{\tilde{n}}$ verifies:

- $\mu_{\tilde{n}} : \mathbb{R} \rightarrow [0, 1]$ is a continuous application;
- $\exists ! n \in \mathbb{R} : \mu_{\tilde{n}}(n) = 1$;
- $\forall x, y \in \mathbb{R}; \forall \lambda \in [0, 1]; \mu_{\tilde{n}}(\lambda x + (1 - \lambda)y) \geq \min(\mu_{\tilde{n}}(x), \mu_{\tilde{n}}(y))$ (fig. 1.5).

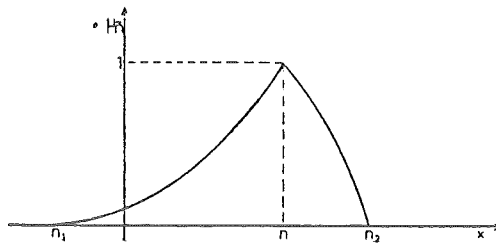


Fig. 1.5 - Membership function of a fuzzy real number

- n is the modal value of \tilde{n} ;
- the fuzzy set $\tilde{n} = \{(x), \mu_{\tilde{n}}(x), x \in \mathbb{R}\}$ represents a real numbers set approximately equal to n and the $\mu_{\tilde{n}}$ form allows the expression of an individual judgement on the reliance level precision attributed to the real x , i.e., $\mu_{\tilde{n}}(x)$ gives the degree according to which "the value of \tilde{n} is x ".
- the last condition in def. 1.13 is equivalent to saying that \tilde{n} is a convex fuzzy set [35].
- we shall designate by $\tilde{\mathcal{N}}(\mathbb{R})$ the set of fuzzy real numbers.

Def. 1.14 - We say that the fuzzy real number, \tilde{n} , is positive (negative) if it's membership function is such that $\mu_{\tilde{n}}(x) = 0$, for all negative or null values x (positive or null).

Then we say that $\tilde{n} \in \tilde{\mathcal{N}}(\mathbb{R}^+)$ ($\tilde{n} \in \tilde{\mathcal{N}}(\mathbb{R}^-)$) and we write $\tilde{n} > 0$ ($\tilde{n} < 0$).

An important field of application of the principle of the extension is in the domain of the algebraic operations with fuzzy real numbers [6, 7, 9].

In general, the binary operation $*$: $\mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ admits a generalization in $\tilde{\mathcal{N}}(\mathbb{R})$, designated by \otimes and established at the cost of the referred principle in the following manner: given \tilde{m}, \tilde{n} ,

$\in \tilde{\mathcal{N}}(\mathbb{R})$, the fuzzy set $\tilde{m} \otimes \tilde{n} = \{(z, \mu_{\tilde{m} \otimes \tilde{n}}(z)), z \in \mathbb{R}\}$ is defined by

$$(1.1) \mu_{\tilde{m} \otimes \tilde{n}}(z) = \sup_{\{(x,y): z=x*y\}} \min(\mu_{\tilde{m}}(x), \mu_{\tilde{n}}(y)), z \in \mathbb{R}$$

In a similar way, if ψ is a unary operation in \mathbb{R} , $\psi: \mathbb{R} \rightarrow \mathbb{R}$, we may consider

$$(1.2) \mu_{\psi \tilde{m}}(z) = \sup_{\{x: z=\psi(x)\}} \mu_{\tilde{m}}(x), z \in \mathbb{R}$$

by application of the expressed principle.

The previous relations permit to define the usual algebraic operations in $\tilde{\mathcal{N}}(\mathbb{R})$ by an adequate particularization of the functions $*$ and ψ .

In this manner, if m and n are real fuzzy numbers, having modal values \underline{m} and \underline{n} respectively, we get:

1. The symmetric of \tilde{n} , designated by $\Theta \tilde{n}$, is obtained by making $\psi(x) = -x$ in (1.2).

Therefore, we conclude that $\Theta \tilde{n}$ is a element of $\tilde{\mathcal{N}}(\mathbb{R})$, since $\mu_{\Theta \tilde{n}}(x) = \mu_{\tilde{n}}(-x), \forall x \in \mathbb{R}$.

2. Considering, for a non null x , $\psi(x) = \frac{1}{x}$ in (1.2), we can characterize the inverse of $\tilde{n}, \tilde{n}^{-1}$. In this case, $\mu_{\tilde{n}^{-1}}(x) = \mu_{\tilde{n}}(\frac{1}{x}), \forall x \neq 0$.

Nevertheless, it can be noticed that $\tilde{n}^{-1} \in \tilde{\mathcal{N}}(\mathbb{R})$ only if $\tilde{n} > 0$ or $\tilde{n} < 0$; otherwise, \tilde{n}^{-1} stops being convex and $\mu_{\tilde{n}^{-1}}$ does not decrease when $|x| \rightarrow \infty$.

3. The multiplication of a scalar by the real fuzzy number \tilde{n} is defined by $\forall x \in \mathbb{R}, \forall \lambda \in \mathbb{R} \setminus \{0\}, \mu_{\lambda \tilde{n}}(x) = \mu_{\tilde{n}}(\frac{x}{\lambda})$, taking $\psi(x) = \lambda x$ in (1.2). Consequently $\lambda \tilde{n} \in$

$\tilde{\mathcal{N}}(\mathbb{R}), \forall \lambda \in \mathbb{R}$ and $\lambda \neq 0$.

4. The defining equality for the addition of \tilde{m} and \tilde{n} is obtained from (1.1) by considering the addition of the reals as particularization of $*$. *Dubois and Prade* [6,7,9] proved

that $\tilde{m} \oplus \tilde{n}$ is also an element of $\tilde{\mathcal{N}}(\mathbb{R})$, associated with the modal value $m + n$. They also established the commutative and associative properties of this operation, having verified that the corresponding neutral element is the real number zero and verified the inexistence of the symmetrical element in the sense of group structure.

One of the consequences of this fact is the difficulty in solving fuzzy equations, usual elimination of terms being "prohibited".

5. The notion of a symmetrical fuzzy number allows to write $\tilde{m} \ominus \tilde{n} = \tilde{m} \oplus (\Theta \tilde{n})$, the membership function of $\tilde{m} \ominus \tilde{n}$ being obtained from (1.1) and the intervening equality understood in the sense of def. 1.5.

Therefore, we conclude that $\tilde{m} \ominus \tilde{n}$ is still a fuzzy real number, this defining the subtraction in $\tilde{\mathcal{N}}(\mathbb{R})$.

6. If, in (1.1) the binary operation $*$ designates the usual product, we get the concept of product of real fuzzy numbers. However, some supplementary conditions are necessary in this case for guaranteeing the stability of the operation. Then:

(i) if $\tilde{m} > 0$ and $\tilde{n} > 0$, then $\tilde{m} \odot \tilde{n}$ is a positive fuzzy number of modal value $m n$ [6,7,9];

(ii) if $\tilde{m} > 0$ and $\tilde{n} < 0$, the equality $\tilde{m} \odot \tilde{n} = \Theta(\tilde{m} \ominus (\Theta \tilde{n}))$ permits to conclude that

$$\tilde{m} \odot \tilde{n} \in \tilde{\mathcal{N}}(\mathbb{R}^-)$$

(iii) if $\tilde{m} > 0$ and $\tilde{n} < 0$, as $\tilde{m} \odot \tilde{n} = (\Theta \tilde{m}) \ominus (\Theta \tilde{n})$, we conclude that $\tilde{m} \odot \tilde{n} \in \tilde{\mathcal{N}}(\mathbb{R}^+)$.

The product of fuzzy numbers satisfies the commutative and associative properties, admitting as neutral element the real 1. Also in this case we recognise the inexistence of the inverse in the sense of group's stability.

7. The division in $\tilde{\mathcal{N}}(\mathbb{R}^+)$ is established from the product operation of fuzzy numbers and from the inverse of a fuzzy real number. As a matter of fact:

(i) if $\tilde{m} > 0$ and $\tilde{n} > 0$, the equality $\tilde{m} \odot \tilde{n} = \tilde{m} \ominus \tilde{n}^{-1}$ considered in the sense of the definition 1.5, shows that $\tilde{m} \odot \tilde{n}$ is still a fuzzy real number. Note that it's first member can be obtained once more from (1.1), when the binary operation $*$ is reduced to the quotient of real numbers.

(ii) In the remaining cases ($\tilde{m} > 0$ and $\tilde{n} < 0$, $\tilde{m} < 0$ and $\tilde{n} > 0$, $\tilde{m} < 0$ and $\tilde{n} < 0$), we may conclude an analogous result, based on the identity $\Theta(\tilde{n}^{-1}) = \Theta(\tilde{n}^{-1})$.

Dubois and Prade [7,9] present a general algorithm for the execution of the generalised operations in $\tilde{\mathcal{N}}(\mathbb{R})$. Having in mind to improve the computational efficiency of these operations in problems of real dimension, the authors [6,7,9]

have also suggested a specific representation of the elements $\tilde{\mathcal{N}}(\mathbb{R})$, that revealed to be indispensable for the algebraic calculation.

Def 1.15 - A real fuzzy number, \tilde{n} , of modal value n , is named of LR type if for all real x

$$\mu_{\tilde{n}}(x) = \begin{cases} L\left(\frac{n-x}{\alpha}\right) & \text{if } x \leq n \\ R\left(\frac{x-n}{\beta}\right) & \text{if } x > n \end{cases}, \alpha, \beta > 0$$

L and R being even non decreasing real functions of real variable, in $[0, +\infty[$ such that $L(0) = R(0) = 1$, left and right reference functions of \tilde{n} , respectively.

Then we write $\tilde{n} = (n, \alpha, \beta)_{LR}$.

Obs 1.6:

(i) α and β are designated left and right limit of \tilde{n} , respectively, and the higher these values the more fuzzy \tilde{n} we get;

(ii) by convention, every real rigid number w admits the representation $w = (w, 0, 0)_{LR}$, L and R being arbitrary.

Based on the previous definition, *Dubois and Prade* [6,7,9] derived rather efficient calculation formulae for the manipulation of the considered fundamental operations.

Theorem 1.1 - Let $n = (n, \alpha, \beta)_{LR}$. Then:

(i) $\Theta n = (\tilde{n}, \beta, \alpha)_{RL}$;

(ii) if $\tilde{n} > 0$ we get $n^{-1} \cong \left(\frac{1}{n}, \frac{\beta}{n^2}, \frac{\alpha}{n^2}\right)_{PL}$, in the proximity of $\frac{1}{n}$;

(iii) $\lambda_{\tilde{n}} = \begin{cases} (\lambda n, \lambda \alpha, \lambda \beta)_{LR} & \text{if } \lambda \geq 0 \\ (\lambda n, -\lambda \beta, -\lambda \alpha)_{RL} & \text{if } \lambda < 0 \end{cases}$

Note that the defining formula of the inverse of a negative fuzzy number is similar to the previous one because $\Theta(\tilde{n}^{-1}) = \Theta(\tilde{n}^{-1})$ both being approximate relations.

Theorem 1.2 - Let $\tilde{m} = (m, \alpha, \beta)_{LR}$, $\tilde{n} = (n, \gamma, \delta)_{LR}$ and $k = (k, \Theta, z)_{RL}$.

Then:

(i) $\tilde{m} \oplus \tilde{n} = (m + n, \alpha + \gamma, \beta + \delta)_{LR}$;

(ii) $\tilde{m} \ominus k = (m - k, \alpha + \zeta, \beta + \Theta)_{LR}$;

(iii) $\tilde{m} \odot \tilde{n} \cong (mn, m\gamma + n\alpha, m\delta + n\beta)_{LR}$ if $\tilde{m} > 0, \tilde{n} > 0$

$k \odot \tilde{n} \cong (kn, n\Theta - k\delta, n\zeta - k\gamma)_{RL}$ if $k < 0, \tilde{n} > 0$

$\tilde{m} \odot \tilde{n} \cong (mn, -n\beta - m\delta, -n\alpha - m\gamma)_{RL}$ if $\tilde{m} < 0, \tilde{n} < 0$

(iv) $\tilde{m} \odot \tilde{n} \cong \left(\frac{m}{k}, \frac{m\zeta + k\alpha}{k^2}, \frac{m\Theta + k\beta}{k^2}\right)_{LP}$, on the proximity of $\frac{m}{k}$, if $\tilde{m} > 0$

and $k > 0$.

Observ 1.7:

(i) The presented addition and subtraction formulae are accurate. In particular, if \tilde{m} and k are of LL or RR type, the calculation of the subtraction is largely simplified.

(ii) As for the calculation of $\tilde{m} \odot \tilde{n}$, the accurate formulae that we have referred are the more corrected the lower the left and right limits are with respect to the modal values m and n . When such limits can not be disregarded, other approximate formulae may be used. [9].

(iii) If division in $\tilde{\mathcal{N}}(\mathbb{R})$ is defined from the notion of inverse of a fuzzy number and from the product operation, the corresponding calculation formula is also an approximation. The same argument justifies that the division \tilde{m} by \tilde{k} , when $\tilde{m} < 0$ or $\tilde{k} < 0$, is translated by expressions similar to the ones considered in the previous theorem obtained from the corresponding formulae for the multiplication.

It can be verified, that LR representation of fuzzy numbers constitutes a powerful instrument for executing the fundamental algebraic operations, since the manipulation of the intervening parameters is sufficient to obtain the final membership function and the established expressions do not depend on the particular analytic form of the L and R functions.

Obs. 1.8:

The present concept of real fuzzy number admits several important generalizations:

- (i) The notion of a fuzzy interval is also due to *Dubois and Prade* [6,7,9]. According to this concept, it is assumed that there exists an interval of modal values, $[\underline{n}, \bar{n}]$, in which the membership function takes constantly the value 1, in contrast with what happens in definition 1.13, where we assume the uniqueness of n 's modal value (fig. 1.6).

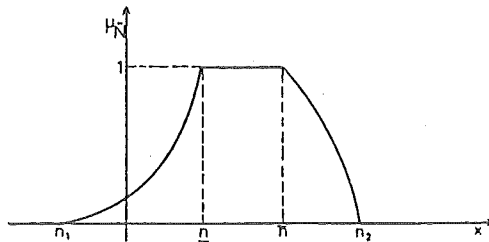


Fig. 1.6 - Membership function of a fuzzy interval

The algebraic operations considered in the previous section are easily generalized to the new entity, that also admits a LR representation, including the concepts of fuzzy real number and rigid real interval.

- (ii) On the other hand, in more recent works [5,25] the hypothesis of continuity in definition 1.13 are relaxed. In effect, that definition seems in same way incomplete as it does not permit to consider the rigid real set numbers as a subset of $\tilde{\mathcal{N}}(\mathbb{R})$.
- (iii) Finally, it comes the type 2 fuzzy number notion [6,7,9] as a particularization of the concept of type 2 fuzzy set, previously referred.

2 - Introduction to Decision Theory in a Fuzzy Environment

In a conventional environment, we usually consider a set of constraints, delimiting the "rigid" set of the admissible decisions and one (or more) functions that order the diverse alternatives according to some criteria, always supposing that the intervenient data is given and perfectly known and that those functions are themselves formalized.

Nevertheless, it is well known that in real situations these assumptions must be relaxed. The approach, based on the fuzzy sets theory gives in this domain an indispensable contribution.

2.1 - The decision concept in a difuse environment

In 1970, *Bellman and Zadeh* suggested a decision model in a difuse environment, possessing, in it's nature, the same structure as the classic model and that was the starting point for a great part of the research in this matter.

X being a universe of possible alternatives (it's elements are objects about which restrictions and objectives of certain optimization problem are formulated). The authors [2] considered a decision model in which the objectives and the restrictions are fuzzy. (A fuzzy objective, \tilde{G} , and a fuzzy restriction, \tilde{C} , are fuzzy sets in X defined by membership functions $\mu_{\tilde{G}} : X \rightarrow [0,1]$ and $\mu_{\tilde{C}} : X \rightarrow [0,1]$ respectively). Once we intend to satisfy certain objectives as well as the constraints, we define decision on a fuzzy environment as a selection of the alternatives that simultaneously satisfy constraints and objectives. Therefore, they correspond the "logic and" to the intersection of the sets, in this way turning the relationship between the restrictions and the objectives a symmetrical one.

Def 2.1 - Let G_j ($j = 1, \dots, n$) be n fuzzy objectives and C_i ($i = 1, \dots, m$) m fuzzy constraints. A fuzzy decision is a fuzzy set in X , defined by:

$$D = \left(\bigcap_{j=1}^m C_i \right) \cap \left(\bigcap_{j=1}^n G_j \right),$$

with the membership function

$$\mu_D(x) = \left(\bigwedge_{i=1}^m \mu_{C_i}(x) \right) \wedge \left(\bigwedge_{j=1}^n \mu_{G_j}(x) \right), \quad \forall x \in X,$$

where $a \wedge b = \min(a,b)$.

Ex 2.1 - The board of directors of an enterprise must decide about the dividend that must be offered to the stockholders on the next meeting. Some members discussed that the dividend must be substantially higher than 10% in order to make attractive profits. We take this requirement as the objective. The membership function of the fuzzy set "dividend substantially higher than 10" can be analytically defined by:

$$\mu_{\tilde{G}}(x) = \begin{cases} 0 & \text{if } x \leq 10 \\ 1/[1+(x-10)^{-2}] & \text{if } x \geq 10 \end{cases} \quad (\text{fig.2.1})$$

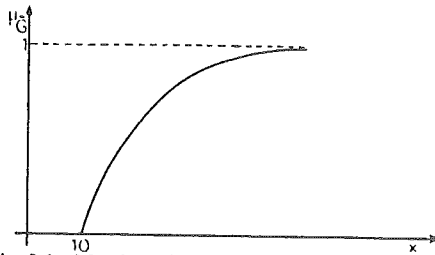


Fig. 2.1 - Membership function of the fuzzy objective

Nevertheless, other members of the board insisted that, by stability reasons, "the dividend must be on the proximity of 11%". This condition can be considered as a fuzzy restriction, whose membership function can be translated, for example by:

$$\mu_{\tilde{C}}(x) = \frac{1}{1+(x-11)^4}, \quad x > 0 \quad (\text{fig.2.2})$$

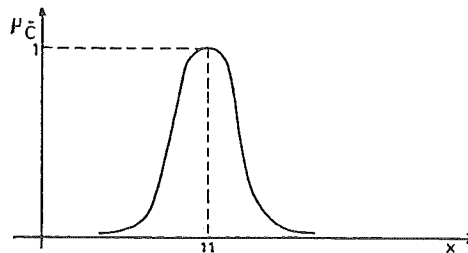


Fig. 2.2 - Membership function of a fuzzy restriction

The membership function of the fuzzy decision is therefore defined by

$$\mu_D(x) = \min(\mu_G(x), \mu_S(x)), \quad x > 0 \quad (\text{fig. 2.3})$$

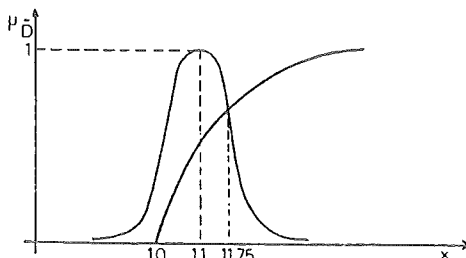


Fig. 2.3 - Membership function of the fuzzy decision

Generally, we opt for a rigid decision. It seems adequate to consider as an optimal rigid decision the solution with the greater membership level relatively to the fuzzy decision, i.e., $x^* \in X$ as well that $\mu_D(x^*) \geq \mu_D(x), \forall x \in X$.

The definition 2.1 is based on three fundamental conditions: the aggregation of objectives and constraints, possessing equal importance corresponds to the "logical and"; the "logical and" is translated by the set intersection and this is defined by the min operator. These positions are not universally accepted.

In the first place, (this has been previously referred), we can consider other operators (Hamacher, algebraic product, etc.) to model the fuzzy sets intersection. On the other hand, it is still questionable if the intersection constitutes a type of adequate aggregation or if we should adopt more general confluence forms.

Ex 2.2 - A professor must decide how to quote the solution of a linear program proposed to his students. To answer the question, the student can apply the graphic method or the simplex method. Two fuzzy sets were defined: "acceptable graphic solution" (G) and "acceptable algebraic solution" (S) and the student's level is characterized by the attribution degree obtained relatively to each of the sets.

Being $\mu_G = 0.9$ and $\mu_S = 0.7$.

If the global value of the question would correspond to the attributed degree in relation to the fuzzy set "acceptable solutions to linear programming" it wouldn't be strange that, for the student in question, this level should be determined by

$$\mu_{LP}(x) = \max(\mu_G(x), \mu_S(x)) = 0.9,$$

i. e., defining the decision at the cost of the max operator.

We observe then that the mathematic representation of the decision concept must be flexible, suffering modifications depending on the context in which it is inserted. In this sense, we can point out the following suggestions:

- (i) *Bellman and Zadeh* proposed the attribution of weights to the membership functions of the objectives and constraints in a decision process, according to its relative importance. Appearing, in this way, the concept of convex decision [2].
- (ii) *Yager*, in an effort of adapting the models of "and" and "or" to the real context in which they are used for, make the min and max operators flexible, presenting new operators to define the reunion and intersection of fuzzy sets [34].
- (iii) *Zimmermann* [40,41] distinguishes then clearly: when we interpret the decision notion as fuzzy set intersection, defined by the min or product operators, we do not admit any compensation between the attribution levels of the sets involved, since any of these operators leads to attributed degrees less than or equal to the minor degree reached by the intersected set (ex 2.1); on the contrary, if the decision is translated by the sets reunion, through the max operators, we are lead to a total compensation of low attributed values by the maximum attribution degree (ex 2.2). Nevertheless, when we consider business decisions, we observe that generally there are compensations, either between the attributed degree of the objectives or in the levels according to which the constraints limit the extent of the

decisions, but those compensations are not total. The author proposed to test empirically a new type of aggregation, denominated "and compensatory", allowing some compensation degree, i.e. standing between "and" not compensatory and the "or totally compensatory. This lead to the y operator.

- (iv) Other authors [18], pointed out the numerical difficulties risen by the y operator and present a new aggregation rule named min-limited sum operator.
- (v) Finally, in light of a more general situation, on which G_j ($j = 1, \dots, n$) are fuzzy objectives for the alternatives set Y and the fuzzy constraints G_i ($i = 1, \dots, m$) are fuzzy sets in X ($X \neq Y$) more general operators can be adopted, leading to new decision models [9, 23].

We turn up with a great variety of operators for making the aggregations of fuzzy sets, fact that can make difficult the choice of certain operator in a particular situation. Zimmermann [41] establishes and discusses 8 useful classifications criteria on the selection of an adequate connective operator.

In a general way, the "confluence of objectives and constraints" (that leads to a decision D) can be formulated by:

$$\mu_D(x) = \Phi (\mu_{G_1}(x), \dots, \mu_{G_n}(x), \mu_{C_1}(x), \dots, \mu_{C_m}(x)), \forall x \in X,$$

where Φ designates a function of an adequate aggregation.

2.2 - Fuzzifying a linear program

The L.P. models can be considered as a particular type of decision model, admitting the following formulation:

$$(2.1) \quad \max c^t x \quad \text{s.a.} \quad Ax \leq b, \quad x \in X,$$

where $c \in \mathbb{R}^n$, $b \in \mathbb{R}^m$ and $A \in \mathbb{R}^{m \times n}$ have the usual meaning and X is a subset of \mathbb{R}^n defined by linear constraints and eventually conditions of type ≥ 0 .

According to the classification made by Negoita [22], we point out two great approaches whose nature can be distinguished by the way they introduce the diffusion in the model:

1. The flexible programming ("soft" on Dubois's and Prade's terminology [9]) admits the inflexibility of the coefficients in the model (b , c and A components), making weak the mathematic relations that take part in them. In this way:
 - (i) What respect to o.f., the weakening of the rigid criterion of optimization is obtained by the relaxation of the "maximum value" concept. The decision maker gets "a highly sufficient value" in the constraints set of the problem, reaching certain level of satisfaction established by him.

Using "-" as a diffusion symbol, we write $\max x^t c$.

- (ii) The fuzzifying of the problem constraints (2.1) is obtained by a flexible relation " \leq ". In a fuzzy environment, the decision maker accepts a few constraint violations, but attributes them different degrees of importance, according to the occurred constraints. In this manner, the rigid condition $a_j^t x \leq b_j$ is substituted by the fuzzy condition " $a_j^t x \leq b_j$ ", whose meaning may be, for example, " $a_j^t x$ essentially less than or equal to b_j ", allowing the set $\{x \in \mathbb{R}^n : a_j^t x \leq b_j\}$ to admit not only the vectors that satisfy the rigid condition $a_j^t x \leq b_j$, but also vectors for which $a_j^t x > b_j$, since that the corresponding violation does not exceed a value established by it, called maximum degree of tolerance. However, not remaining indifferent to the vectors that are on the 1st and 2nd cases, it expresses its preference through a real function that translates of the degree to which a certain vector $x \in \mathbb{R}^n$ satisfies the constraint $a_j^t x \leq b_j$. In this manner, this fuzzy constraint is mathematically represented by the fuzzy set \tilde{C}_j , and it can be said that $\{x \in \mathbb{R}^n : a_j^t x \leq b_j\} = \{(x, \mu_{\tilde{C}_j}(x)), x \in \mathbb{R}^n\}$.

- 2 - The robust programming also allow to model problems whose structures are not exactly known, but takes into account inaccuracy of the data, immediately on the building stage of the model. The coefficients of the model are considered fuzzy parameters, which originate the following formulation:

$$\max c^t x \quad \text{s.a.} \quad \tilde{A}x \leq b, \quad x \in X.$$

It is however, necessary to redefine the mathematical language used to express the subjacent model (namely the meaning of the relationship $\leq, =, \geq$).

Obs 2.1 :

- (i) A more general characterization abandoned the punctual character of the searched $x \in X$ solution, attributing to it a fuzzy character. The solution to the corresponding fuzzy L.P. is, in this case, a fuzzy set in R^n , the symbol of the difusion being placed under the variable of the problem, $x \in \tilde{P}(R^n)$;
- (ii) on the two approaches we should decide upon the function types used every time we want to represent, in the context of the fuzzy sets theory, quantities or relations that are not known with precision.
- (iii) We adopted the same notation to represent the "minimization" relaxed operation, as well as the other constraints types, writing respectively $\min c^t x, a_i^t x \geq b_i, a_i^t x \leq b_i$.

Therefore we have noticed that on opposition to the classic L.P., the fuzzy L.P. does not constitute a uniquely defined model, allowing a reasonable number of semantic variations formally dependent on the hypothesis and characteristics of the real situation being modeled.

3 - Flexible Fuzzy L.P.

It is the approach more close to the model proposed by *Bellman* and *Zadeh* once that, in general, it makes the objective function and/or the constraints be represented by fuzzy sets, whose aggregation leads to the fuzzy decision concept.

We distinguish two approaches from now on: the model totally fuzzy (on which the objective and the constraints are fuzzy) and the model partially fuzzy (for which one intends to optimize one o.f. classic in a domain defined by a set of fuzzy constraints).

3.1 - The totally fuzzy model

L.P. being fuzzy:

$$\begin{aligned}
 (3.1) \quad \max \quad c^t x \quad \text{s.a.} \quad & a_i^t x \geq b_i \quad i = 1, \dots, k_1 \\
 & a_i^t x \leq b_i \quad i = k_1 + 1, \dots, k_2 \\
 & a_i^t x = b_i \quad i = k_2 + 1, \dots, k_3 \\
 & x \in X,
 \end{aligned}$$

where X is the rigid subset of R^n already mentioned.

We represent, as it has been already said, the i th fuzzy constraint by a fuzzy set in R^n , \tilde{C}_i ($i = 1, \dots, k_3$). Generally this set is normalized, being the corresponding membership function $\mu_{\tilde{C}_i} : R^n \rightarrow [0,1]$ defined, for all x , by $\mu_{\tilde{C}_i}(x) = h(p_i, b_i)(a_i^t x)$, in what $h(p_i, b_i) : R \rightarrow [0,1]$ is a continuous application whose expression depends on the type of constraint in question. The parameter p_i which is associated with it is a positive real number translating of the maximal level of the corresponding tolerance. For example, for one constraint of the type $a_i^t x \leq b_i$, $\mu_{\tilde{C}_i}$ could take the value "1" if the constraint was totally satisfied ($a_i^t \leq b_i$); "0" in the case of being strongly violated ($a_i^t > b_i + p_i$) and decreasing monotonically in the interval $]b_i, b_i + p_i]$.

Analogously, we can represent the o.f. by the fuzzy set $\tilde{C} = \{(x, \mu_{\tilde{C}}(x)), x \in R^n\}$, calibrated by an aspiration level z_0 ($z_0 \in R$) established by the decision maker, by writing according to the considerations made on the preceeding section:

$$(3.2) \quad \max_x c^t x \quad \text{is equivalent to determine } x \text{ such as } c^t x \geq Z_0.$$

Designating by \tilde{D} the fuzzy set defined due to an aggregation function Φ , we've got:

$$\tilde{D} = \{(x, \mu_{\tilde{D}}(x)), x \in R^n\}, \text{ with } \mu_{\tilde{D}} = \Phi(\mu_{\tilde{C}}, \mu_{\tilde{C}_1}, \dots, \mu_{\tilde{C}_{k_3}}).$$

Def 3.1 - We say that $x_0 \in R^n$ is a punctual solution of (3.1) if $x_0 \in X$ and $\mu_{\tilde{D}}(x_0) \geq \mu_{\tilde{D}}(x)$, $\forall x \in X$.

Therefore, to find the punctual solution of (3.1), we will solve the following optimization problem:

$$(3.3) \quad \max_{x \in X} \mu_{\tilde{D}}(x),$$

whose complexity depends fundamentally on the adopted Φ operator and on the type of membership functions considered in the representation of the intervenient fuzzy sets.

The most frequent approach accepts the congruence of the "logical and" with the intersection of the sets and represents it by the min operator; in this version, the problem (3.3) leads to the particular form:

$$(3.4) \quad \max_{x \in X} \mu_G(x) \wedge \left(\bigwedge_{i=1}^{k_3} \mu_{C_i}(x) \right)$$

The solution of (3.4) can be achieved by diverse approaches, in which we can point out the ones of *H. J. Zimmermann and Tanaka, Okuda and Asai*.

3.1.1. - The symmetrical model

Zimmermann [33, 38, 39, 41] based on the transformation operated in (3.2), considers the o.f. as a fuzzy constraint, establishing a symmetrical model, in which the objectives and constraints can become indistinguishable. According to this perspective, we propose to solve the following system of fuzzy constraints:

$$(3.5) \quad \text{determine } x \text{ such as } \begin{aligned} a_i^1 x &\geq b_i & i = 0, \dots, k_1 \\ a_i^1 x &\leq b_i & i = k_1 + 1, \dots, k_2 \\ a_i^1 x &\cong b_i & i = k_2 + 1, \dots, k_3 \\ x &\in X, \end{aligned}$$

where $a_0 = c$ and $b_0 = z_0$.

On these conditions, the problem (3.4) reduces to:

$$(3.6) \quad \max_{x \in X} \bigwedge_{i=0}^{k_3} \mu_j(x), \text{ being } \mu_0 \equiv \mu_G \text{ and } \mu_j \equiv \mu_{C_i}, 1 \geq 1.$$

One of the fundamental transformations of (3.6) is obtained from

Theorem 3.1 - x^* is the optimal solution of (3.6) iff (x^*, λ^*) , with $\lambda^* = \bigwedge_{i=0}^{k_3} \mu_i(x^*)$ is the optimal solution to the program:

$$(3.7) \quad \begin{aligned} &\max \lambda \\ \text{s.a.} \quad &\lambda \leq \mu_j(x) \quad i = 0, \dots, k_3 \\ &x \in X, \lambda \in [0, 1]. \end{aligned}$$

The structure of the previous problem (linearity, convexity,...) depends essentially on the type of the adopted membership functions. The more usual models [17, 33, 38, 39, 41] consider linear membership functions.

Let:

$$(3.8) \quad \mu_i(x) = h^1_{(p_i, b_i)}(a_i^1 x) = \begin{cases} 1 & \text{if } a_i^1 x \geq b_i \\ 1 + (a_i^1 x - b_i) / p_i & \text{if } b_i - p_i \leq a_i^1 x < b_i \\ 0 & \text{if } a_i^1 x < b_i - p_i \end{cases} \quad (\text{fig.3.1})$$

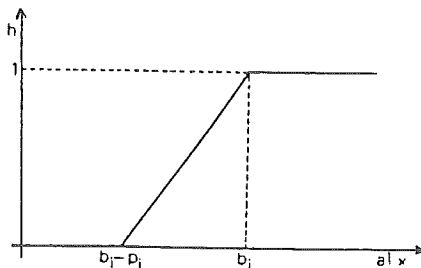


Fig. 3.1 - Linear membership function of the restriction $a_i^1 x \geq b_i$

$$(3.9) \quad \mu_i(x) = h^2_{(p_i, b_i)}(a_i^l x) = \begin{cases} 1 & \text{if } a_i^l x \leq b_i \\ 1 - (a_i^l x - b_i) / p_i & \text{if } b_i \leq a_i^l x \leq b_i + p_i \\ 0 & \text{if } a_i^l x > b_i + p_i \end{cases} \quad (\text{fig.3.2})$$

and

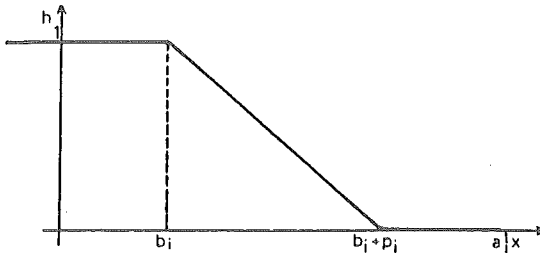


Fig. 3.2 - Linear membership function of the restriction $a_i^l x \leq b_i$

$$(3.10) \quad \mu_i(x) = h^3_{(p_i, \bar{p}_i, b_i)}(a_i^l x) = \begin{cases} 0 & \text{if } a_i^l x < b_i - p_i \\ 1 + (a_i^l x - b_i) / p_i & \text{if } b_i - p_i \leq a_i^l x \leq b_i \\ 1 - (a_i^l x - b_i) / \bar{p}_i & \text{if } b_i < a_i^l x \leq b_i + \bar{p}_i \\ 0 & \text{if } a_i^l x > b_i + \bar{p}_i \end{cases} \quad (\text{fig.3.3})$$

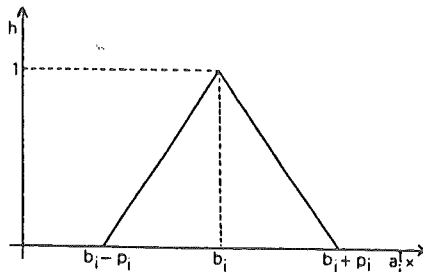


Fig. 3.3 - Linear membership function of the restriction $a_i^l x \cong b_i$

be the membership functions associated respectively with the constraints $a_i^l x \geq b_i$ ($i = 0, \dots, k_1$), $a_i^l x \leq b_i$ ($i = k_1 + 1, \dots, k_2$) and $a_i^l x \cong b_i$ ($i = k_2 + 1, \dots, k_3$).

Making some elementary transformations and applying the theorem 3.1, we can conclude that the following linear program is equivalent to (3.6):

$$(3.11) \quad \begin{aligned} & \max \lambda \\ \text{s.a.} \quad & \lambda p_i - a_i^l(x) \leq p_i - b_i \quad i = 0, \dots, k_1 \\ & \lambda p_i + a_i^l(x) \leq p_i + b_i \quad i = k_1 + 1, \dots, k_2 \\ & \lambda \bar{p}_i - a_i^l(x) \leq \bar{p}_i - b_i \quad i = k_2 + 1, \dots, k_3 \\ & \lambda \bar{p}_i + a_i^l(x) \leq \bar{p}_i + b_i \quad i = k_2 + 1, \dots, k_3 \\ & x \in X, \quad \lambda \in [0, 1]. \end{aligned}$$

The optimum value of λ constitutes a satisfaction measure of the fuzzy system (3.6), once that if

$$(x^*, \lambda^*) \text{ is the optimal solution to (3.11), therefore } \lambda^* = \bigwedge_{i=0}^{k_3} \mu_i(x^*) = \max_{x \in X} \bigwedge_{i=0}^{k_3} \mu_i(x),$$

indicates the attributed level of the maximizing decision x^* relative to the given system.

Ex 3.1 - Let us consider the system of fuzzy linear constraints

$$(3.12) \quad \begin{aligned} 2x_1 - x_2 &\geq 3 \\ 2x_1 - 7x_2 &\leq 0 \\ x_1 + x_2 &\cong 1 \\ x &\in X, \end{aligned}$$

where $X = \{(x_1, x_2) \in \mathbb{R}^2 : x_1 \geq 0, x_2 \geq 0\}$.

Let us take $p_1 = p_2 = 1, p_3 = 0.5, \bar{p}_3 = 1$ as the maximum levels of tolerance for the 1st, 2nd, 3rd constraints, associating them the membership functions μ_1, μ_2 and μ_3 defined by (3.8), (3.9) and (3.10) respectively.

As it was already said, the solution of a system of this type is obtained by solving the linear program (3.11). Therefore the solution of (3.12) is obtained from

$$(3.13) \quad \begin{aligned} &\max \lambda \\ \text{s.a.} \quad &\lambda - 2x_1 + x_2 \leq -2 \\ &\lambda + 2x_1 + 7x_2 \leq 1 \\ &\lambda - 2x_1 - 2x_2 \leq -1 \\ &\lambda + x_1 + x_2 \leq 2 \\ &x_1, x_2 \geq 0, \lambda \in [0, 1]. \end{aligned}$$

The searched solution is $x^* = (4/3, 7/24)$, being $\lambda^* = \bigwedge_{i=1}^3 \mu_i(x^*) = 3/8$ its attribution level relative to the fuzzy set defined by the constraints of the system in cause.

Obs 3.1 :

- 1 - From (3.11) we can distinguish by two others formulations: *Llena* [17] follows the same approach taking $X = \{x \in \mathbb{R}^n : Dx \leq d, x \geq 0\}$ as a rigid system of constraints and *OhEigearthaigh* [24] prefers rather to consider an unsatisfaction measure, $y = 1 - \lambda$, adopting a formulation equivalent to (3.11).
- 2 - Arguing that, once the optimal solution of (3.11) is determined it is not possible to know directly the particular violation level of the i^{th} constraint for the solution x^* , some authors [13] present a modified version from the presented model.

Let us associate an additional variable $t_i = \max(0, b_i - a_i^1 x)$, respectively $t_i = \max(0, a_i^1 x - b_i)$, with each constraint of the type $a_i^1 \geq b_i$, respectively $a_i^1 x \leq b_i$, indicating the corresponding violation level. In relation to a constrain of type $a_i^1 x \cong b_i$, we consider simultaneously two variables $t_i^+ = \max(0, a_i^1 x - b_i)$ and $t_i^- = \max(0, b_i - a_i^1 x)$ that translate, in alternative, the violation level to the right and to the left of b_i , respectively.

Let us define the sets:

$$S_{b-p} = \bigcap_{i=0}^{k_1} \{x \in \mathbb{R}^n : a_i^1 x \geq b_i - p_i\}, S_{b+p} = \bigcap_{i=k_1+1}^{k_2} \{x \in \mathbb{R}^n : a_i^1 x \leq b_i + p_i\}$$

and

$$S_{b \pm p} = \bigcap_{i=k_2+1}^{k_3} \{x \in \mathbb{R}^n : b_i - p_i \leq a_i^1 x \leq b_i + p_i\}$$

If we take into the expression of $\mu_i(x)$, to $i = 0, \dots, k_3$, we can write (3.6):

$$(3.14) \quad \max_{x \in X} \bigwedge_{i=0}^{k_3} \mu_i(x) = \max_{x \in X \cap S_{b-p} \cap S_{b+p} \cap S_{b \pm p}} \bigwedge_{i=0}^{k_3} \mu_i(x)$$

We note that

$$\begin{aligned} \forall x \in S_{b-p}, \quad \mu_i(x) &= 1 - [\max(0, b_i - a_i^{\downarrow}x)/p_i] && \text{for } i = 0, \dots, k_1, \\ \forall x \in S_{b+p}, \quad \mu_i(x) &= 1 - [\max(0, b_i - a_i^{\downarrow}x)/p_i] && \text{for } i = k_1+1, \dots, k_2, \\ \forall x \in S_{b\pm p}, \quad \mu_i(x) &= [1 - (\max(0, b_i - a_i^{\downarrow}x)/p_i)] \wedge [1 - (\max(0, b_i - a_i^{\downarrow}x)/\bar{p}_i)] && \text{for } i = k_2+1, \dots, k_3 \end{aligned}$$

Introducing the defined variables t_i, t_i^{\dagger} and \bar{t}_i , we get from (3.14):

$$\begin{aligned} \max \quad & \left[\bigwedge_{i=0}^{k_2} (1 - t_i/p_i) \right] \wedge \left[\bigwedge_{i=k_2+1}^{k_3} (1 - \bar{t}_i/\bar{p}_i) \right] \wedge \left[\bigwedge_{i=k_2+1}^{k_3} (1 - t_i^{\dagger}/\bar{p}_i) \right] \\ \text{s.a.} \quad & a_i^{\downarrow}x + t_j \geq b_i \quad i = 0, \dots, k_1 \\ & a_i^{\downarrow}x - t_j \leq b_i \quad i = k_1+1, \dots, k_2 \\ & a_i^{\downarrow}x - b_j = t_i^{\dagger} - \bar{t}_i \quad i = k_2+1, \dots, k_3 \\ & 0 \leq t_i \leq p_i \quad i = 0, \dots, k_2 \\ & 0 \leq \bar{t}_i \leq \bar{p}_i \quad i = k_2+1, \dots, k_3 \\ & 0 \leq t_i^{\dagger} \leq \bar{p}_i \quad i = k_2+1, \dots, k_3 \\ & t_i^{\dagger} \cdot \bar{t}_i = 0 \quad i = k_2+1, \dots, k_3 \end{aligned}$$

Finally, by application of theorem 3.1. we get to the problem:

$$\begin{aligned} (3.15) \quad & \max \lambda \\ \text{s.a.} \quad & \lambda p_i + t_i \leq p_i \quad i = 0, \dots, k_2 \\ & \lambda p_i + \bar{t}_i \leq \bar{p}_i \quad i = k_2+1, \dots, k_3 \\ & \lambda p_i + t_i^{\dagger} \leq \bar{p}_i \quad i = k_2+1, \dots, k_3 \\ & a_i^{\downarrow}x + t_i \geq b_i \quad i = 0, \dots, k_1 \\ & a_i^{\downarrow}x - t_i \leq b_i \quad i = k_2+1, \dots, k_2 \\ & a_i^{\downarrow}x - b_i = t_i^{\dagger} - \bar{t}_i \quad i = k_2+1, \dots, k_3 \\ & 0 \leq t_i \leq p_i \quad i = 0, \dots, k_2 \\ & 0 \leq \bar{t}_i \leq \bar{p}_i \quad i = k_2+1, \dots, k_3 \\ & 0 \leq t_i^{\dagger} \leq \bar{p}_i \quad i = k_2+1, \dots, k_3 \\ (3.16) \quad & t_i^{\dagger} \cdot \bar{t}_i = 0 \quad i = k_2+1, \dots, k_3 \\ & x \in X, \quad \lambda \in [0, 1]. \end{aligned}$$

where the condition $x \in X \cap S_{b-p} \cap S_{b+p} \cap S_{b\pm p}$ was omitted as it became redundant.

Let $(x^*, \lambda^*, t, t^{\dagger}, \bar{t})$ be an optimal solution to (3.15) that can or can't satisfying (3.16). Then we know that:

$$\lambda^* = \left[\bigwedge_{i=0}^{k_2} (1 - t_i/p_i) \right] \wedge \left[\bigwedge_{i=k_2+1}^{k_3} (1 - \bar{t}_i/\bar{p}_i) \right] \wedge \left[\bigwedge_{i=k_2+1}^{k_3} (1 - t_i^{\dagger}/\bar{p}_i) \right],$$

where $t_i \geq \max(0, b_i - a_i^{\downarrow}x^*)$ for $i = 0, \dots, k_1$, $t_i \geq \max(0, a_i^{\downarrow}x^* - b_i)$ for $i = k_1+1, \dots, k_2$, $\bar{t}_i \geq \max(0, b_i - a_i^{\downarrow}x^*)$ for $i = k_2+1, \dots, k_3$ and $t_i^{\dagger} \geq \max(0, a_i^{\downarrow}x^* - b_i)$ for $i = k_2+1, \dots, k_3$.

We define $\hat{t}_i = \max(0, b_i - a_i^{\downarrow}x^*)$, for $i = 0, \dots, k_1$, $\hat{t}_i = \max(0, a_i^{\downarrow}x^* - b_i)$ for $i = k_2+1, \dots, k_3$ and

$$\hat{\lambda} = \left[\bigwedge_{i=0}^{k_2} (1 - \hat{t}_i/p_i) \right] \wedge \left[\bigwedge_{i=k_2+1}^{k_3} (1 - \hat{t}_i^{\dagger}/\bar{p}_i) \right] \wedge \left[\bigwedge_{i=k_2+1}^{k_3} (1 - \hat{t}_i/\bar{p}_i) \right].$$

On these conditions, the solution $(x^*, \hat{\lambda}, \hat{t}, \hat{t}^{\dagger}, \hat{t})$ verifies $\hat{\lambda} \geq \lambda^*$ and is admissible for (3.15), therefore satisfying $t_i^{\dagger} \cdot \bar{t}_i = 0$ for $i = k_2+1, \dots, k_3$. This vector is then an optimal solution of (3.15) which satisfies (3.16) implicitly.

In this way, once the solution to problem (3.15) has been determined with t_i, t_i^{\dagger} and \bar{t}_i defined as explained, we know immediately the violation level that for each constraint of the fuzzy system data.

The obtained presents greater dimension (3.11), in spite of the constraints $t_i \leq p_i (i=0, \dots, k_2)$ and $t_i^+ \leq \bar{p}_i, t_i^- \leq \underline{p}_i (i = k_2+1, \dots, k_3)$ becoming redundant.

- 3 - Based on this approximation, some authors studied the structure of the linear programs achieved, namely conditions of existence and optimality [17], as well as the sensitivity of the solution of the fuzzy system to small little variations on the maximum levels of tolerance [13, 17].
- 4 - On the attempt of getting a better approximation of the mathematical models to the real situation of uncertainty, some authors have introduced non-linear membership functions, whose interest increases in cases where it keeps the linearity of the resulting problems.

(i) In the domain of multicriteria optimization and based on the previous empiric results *Leberling* [16] presents an hyperbolical function defined, for all $x \in \mathbb{R}^n$, by:

$$\mu_i(x) = h(\alpha_i, k_i)(a_i^1 x) = \frac{1}{2} \frac{e^{(a_i^1 x - k_i)\alpha_i} - e^{-(a_i^1 x - k_i)\alpha_i}}{(a_i^1 x - k_i)\alpha_i - (a_i^1 x - k_i)\alpha_i} + \frac{1}{2},$$

where α_i is a real parameter and k_i a constant such that $h(k_i) = 0.5$.

This type of non linear functions can be used to define a membership function of a fuzzy set representative of an inequality constraint, taking $a_i > 0$ and $k_i = b_i - p_i/2$ ($a_i < 0$ and $k_i = b_i + p_i/2$) we associate them with a constraint of type $a_i^1 \geq b_i$ ($a_i^1 \geq b_i$).

In this manner, if only constraints of inequality type are present in the model, we still can get a linear program, by application of theorem 3.1, by making the change of the variable $x_{n+1} = \tanh^{-1}(2\lambda - 1)$.

(ii) For representing a fuzzy constraint to the equality type, $a_i^1 \equiv b_i$, *Llena* [17] proposed two types of non linear functions.

- $\mu_i^1(x) = e^{-k_i | a_i^1 x - b_i |}$

and

- $\mu_i^2(x) = [1 + k_i(a_i^1 - b_i)^2]^{-1}, x \in \mathbb{R}^n$,

k_i being a real positive constraint. Still in this case and through an adequate changing of variable, we come up with a linear model.

(iii) An intermediate extension was accomplished by *Nakamura* [20] that associated a stepwise linear function with each fuzzy constraint. The author shows how the solution of the fuzzy system is obtained by solving a sequence of linear programs.

3.1.2. - The assymmetrical model

Tanaka, Okuda and Asai [30] propose a more general model, that emphasises the o.f. and whose formulation is similar to (3.4). They propose to solve the problem:

$$(3.17) \quad \sup_{x \in X} \mu_D(x) = \sup_{x \in X} \mu_G(x) \wedge \mu_C(x), \text{ in which } \mu_C(x) = \bigwedge_{i=1}^{k_3} \mu_{C_i}(x).$$

Using the notion of the α level cut's of a fuzzy set they establish the following fundamental result.

Theor 3.2 - if $X \cap C_\alpha \neq \emptyset$, for all $\alpha \in [0, 1]$, we get:

$$\sup_{x \in X} \mu_G(x) \wedge \mu_C(x) = \sup_{\alpha \in [0,1]} (\alpha \wedge \sup_{x \in X \cap C_\alpha} \mu_G(x)), \text{ with } C_\alpha = \bigcap_{i=1}^{k_3} C_i^\alpha.$$

This property reduces the analysis of (3.17) to the study of a new problem:

$$(3.18) \quad \sup_{\alpha \in [0,1]} (\alpha \wedge \sup_{x \in X \cap C_\alpha} \mu_G(x)).$$

F being $[0, 1] \rightarrow [0, 1]$ defined by $F(\alpha) = \sup_{x \in X \cap C_\alpha} \mu_G(x)$. If we admit the continuity of

F[23], the Brouwer's Theorem of the fixed point allows to conclude the existence of $\alpha \in [0, 1]$ the only value such that $F(\alpha) = \alpha$, where in this case:

$$\sup_{x \in X} \mu_D(x) = \bar{\alpha} = \sup_{x \in X \cap C_{\bar{\alpha}}} \mu_G(x),$$

once,

$$\sup_{\alpha \in [0,1]} (\alpha \wedge F(\alpha)) = \bar{\alpha} = F(\bar{\alpha}).$$

The achieved results do neither depend on the type of constraints and o.f. considered (it can be suitable for more general optimization problems) nor from the type of membership functions adopted. Nevertheless, they reduce the given problem to the determination of α , that is not achieved immediately. On the other hand, they are based on the fundamental hypothesis of the continuity of the function F , that rarely is confirmed in practice.

To answer to the first question, *Negoita and Ralescu* [23] and *Tanaka et al* [30] present formulations equivalent to (3.17), that are structurally more simple. For solving the second question, we point out one sufficient condition of continuity for F , established by *Tanaka et al* [30], and based on the notion of a convex fuzzy set [35].

These authors [30] have developed a complet study of the problem (3.18). For that purpose, they formulated a vast set of hypothesis, that revealed in practice to be very restrictive.

Flachs e Pollatschek [12] presented conditions of existence of the solution to the problem in question, leading to an analysis free of the assumptions of *Tanaka et al*, namely with respect to the continuity of F .

The described approach is illustrated with the following example, involving a linear program and where the membership functions are taken as linear, fact that allows to simplify substantially the discussed questions.

Ex 3.2 - Let the fuzzy linear program be:

$$(3.19) \quad \begin{aligned} \max \quad & Z = 2x_1 + x_2 \\ \text{s.a.} \quad & x_1 + x_2 \leq 4 \\ & x_1 \leq 3 \\ & x_1 + 2x_2 \equiv 6 \\ & x_1, x_2 \geq 0. \end{aligned}$$

Let $z_0 = 14$ be the aspirations level associated with the o.f. and $p_0 = 1, p_1 = 4, p_2 = 6, p_3 = 4$ and $p_3 = 1$ be the maximal tolerance levels associated with the to o.f. and constraints, respectively.

Adopting linear membership functions, we represent the o.f. by the set \tilde{C} , with $\mu_{\tilde{C}}$ given by (3.8) and the constraints by the sets \tilde{C}_1, \tilde{C}_2 and \tilde{C}_3 , with $\mu_{\tilde{C}_1}$ and $\mu_{\tilde{C}_2}$ given by (3.9) and given by (3.10). Let $X = \{(x_1, x_2) \in \mathbb{R}^2 : x_1, x_2 \geq 0\}$ and $\mu_{\tilde{C}}(x) = \bigwedge_{i=1}^{k_3} \mu_{\tilde{C}_i}(x)$, for all the $x \in \mathbb{R}^2$.

According with what we said before, and being $X \cap C_\alpha \neq \emptyset, \forall \alpha \in [0, 1]$, we must determinate the fixed point of F .

Now, by definition of the α level cut, we get, for $\alpha > 0$:

$$C_\alpha = \bigcap_{i=1}^2 \{x \in \mathbb{R}^2 : a_i^1 \leq b_i + p_i(1-\alpha)\} \cap \{x \in \mathbb{R}^2 : b_3 - p_3(1-\alpha) \leq a_3^1 x \leq b_3 + p_3(1-\alpha)\}$$

and therefore F can be obtained solving the L.P. parametric problem:

$$(3.20) \quad \begin{aligned} \max \quad & Z = 2x_1 + x_2 \\ \text{s.a.} \quad & x_1 + x_2 \leq 8 - 4\alpha \\ & x_1 \leq 9 - 6\alpha \\ & x_1 + 2x_2 \leq 10 - 4\alpha \\ & x_1 + 2x_2 \geq 5 + \alpha \\ & x_1, x_2 \geq 0, \alpha \in [0, 1] \end{aligned}$$

whose optimal solution is given by:

$$(3.21) \quad x^*(\alpha) = \begin{cases} (8-4\alpha, 0) & \alpha \leq 1/2 \\ (9-6\alpha, 2\alpha-1) & 1/2 < \alpha \leq 2/3, \text{ with } z^*(\alpha) \\ (11-9\alpha, -3+5\alpha) & 2/3 < \alpha \leq 1 \end{cases} \quad \begin{cases} 16-8\alpha & \alpha \leq 1/2 \\ 17-10\alpha & 1/2 < \alpha \leq 2/3 \\ 19-13\alpha & 2/3 < \alpha \leq 1 \end{cases}$$

Nothing that:

- (i) $F(0) = \mu_{\tilde{C}}(x^*(0)) = 1$;
- (ii) If $x^*(\alpha)$ is the optimal finitr solution of (3.20) and $z^*(\alpha)$ the corresponding optimal value of the o.f., so, for all $\alpha > 0$, and because of $h(z_0, p_0)$ is an increasing

function, we have $\mu_G(x^*(\alpha)) = h(z_0, p_0) (z^*(\alpha)) \geq \mu_G(x^*(\alpha))$, for all admissible solution $x(\alpha)$ of (3.20).

We conclude, for $\alpha \in [0, 1]$, that $F(\alpha) = h(z_0, p_0) (z^*(\alpha))$. F is continuous in $[0, 1]$ because h is continuous in \mathbb{R} and z^* is continuous in $[0, 1]$. It's analytic expression is easily calculated from (3.21):

$$F(\alpha) = \mu_G(x^*(\alpha)) = \begin{cases} 1 & \text{if } \alpha \leq 0.25 \\ 5/3 - 8/3\alpha & \text{if } 0.25 < \alpha \leq 0.5 \\ 2 - 10/3\alpha & \text{if } 0.5 < \alpha \leq 0.6 \\ 0 & \text{if } \alpha > 0.6 \end{cases}$$

Consequently $\bar{\alpha} = 5/11$ is the value that solves the given problem. As solution vector, we take $x^*(\bar{\alpha}) = (68/11, 0)$.

Obs 3.2

- (i) If the original problem is of the minimization type, the same happens with the associated parametric linear program, h being in this case a decreasing function. A similar reasoning lead to similar results.
- (ii) If the associated linear program is unbounded for some value of α , this will occur in the whole interval $[0, 1]$, since, with the formulated hypothesis, it is admissible in the interval. In these conditions, for all $\alpha \in [0, 1]$, there is $x \in X \cap C_\alpha$ such that $c^t x \geq z_0$, i.e., there is $x \in X \cap C_\alpha$ such that $\mu_G(x) = 1$. So $F(\alpha) = 1$, $\forall \alpha \in [0, 1]$ and the problem is solved trivially since $\bar{\alpha} = 1$.

3.2. - The partially fuzzy model

Let us consider a model in which the o.f. is rigid, in the sense that is effectively maximized or minimized, but where the constraints are totally or partially fuzzy:

$$(3.22) \quad \sup_{x \in \bar{C}} f(x)$$

where $X \subset \mathbb{R}^n$, $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and \bar{C} is a fuzzy set in \mathbb{R}^n , defined by $\mu_{\bar{C}}: \mathbb{R}^n \rightarrow [0, 1]$.

The symmetry referred in the previous sections ceases to make sense, since the fuzzy constraints define the decision space and the o.f. induces to an order of the alternatives (as the same as in the classic L.P.). In this manner, it isn't directly applicable the reasoning developed for obtaining the previous formulations.

To approach this type of problems, we can distinguished two approximations, according to the classification given by Zimmermann [41].

3.2.1. - Linear Program with fuzzy solution

Some authors [31] suggested the calculation of the α level cuts on the solution space of the given problem, determinating, for each cut, the corresponding optimal value of the o.f.. They have defined a fuzzy set "optimal values of the o.f.", attributing to each obtained value an attribution degree equal to the associated level cut in the solution space.

We indicate by $X_\alpha (X_\alpha = X \cap C_\alpha)$ the α level cut of the solution space of (3.22).

Def 3.2 - The fuzzy solution of (3.22) is a fuzzy set in \mathbb{R}^n , $\bar{\chi} = \{(x, \chi(x)), x \in \mathbb{R}^n\}$, whose membership function is defined by:

$$(3.23) \quad \chi(x) = \begin{cases} \sup_{\alpha} & \text{if } x \in \cap S(\alpha) \\ \chi \in S(\alpha) & \alpha > 0 \\ 0 & \text{otherwise,} \end{cases}$$

with $S(\alpha) = \{x \in X_\alpha \mid f(x) = \sup_{x' \in X_\alpha} f(x')\}$, $\forall \alpha \in]0, 1]$.

$$\begin{aligned} \max \quad & Z(x_1, \dots, x_n) = \sum_{j=1}^n c_j x_j \\ \text{s.t.} \quad & \sum_{j=1}^n a_{ij} x_j = b_i \quad i = 1, \dots, m \\ & x_1, \dots, x_n \geq 0 \end{aligned}$$

Admitting that the coefficients c_j ($j = 1, \dots, n$), b_i ($i = 1, \dots, m$) and a_{ij} ($i = 1, \dots, m, j = 1, \dots, n$) are $\tilde{N}(\mathbb{R})$ elements, several questions may be raised, regarding the meaning of the existing relations and expressions in the resulting formulation, namely the specification of

$$\sum_{j=1}^n a_{ij} \cdot x_j \quad \text{and} \quad \sum_{j=1}^n c_j \cdot x_j.$$

when x_j ($j = 1, \dots, n$) are rigid or fuzzy variables, the establishment of comparison criteria between fuzzy numbers; and in particular the generalization of the relations " \leq ", " \geq " and " $=$ " and the subsequent adaptation of the objective function.

4.1. Fuzzifying a real linear function

Let $y_i = a_i^t x = \sum_{j=1}^n a_{ij} x_j$ be a linear function in the variables x_1, \dots, x_n .

Suppose that the a_{ij} are fuzzy parameters, mathematically represented by real fuzzy numbers of the same L.R. type, i.e.,

$$\tilde{a}_{ij} = (a_{ij}, \underline{a}_{ij}, \bar{a}_{ij})_{LR}$$

with

$$\underline{a}_{ij}, \bar{a}_{ij} \geq 0 \quad \text{for } j = 1, \dots, n, i = 1, \dots, m$$

and let us search for \tilde{y}_i , the fuzzy version of the considered linear function.

Let x_1, \dots, x_n be non negative scalars. Then, from the definition of scalar product (theorem 1.1) and addition of real fuzzy numbers of LR type (theorem 1.2), we may write:

$$\begin{aligned} \tilde{y}_i &= \tilde{a}_{i1} x_1 \oplus \dots \oplus \tilde{a}_{in} x_n \\ &= (a_{i1} x_1, \underline{a}_{i1} x_1, \bar{a}_{i1} x_1)_{LR} \oplus \dots \oplus (a_{in} x_n, \underline{a}_{in} x_n, \bar{a}_{in} x_n)_{LR} \\ &= \left(\sum_{j=1}^n a_{ij} \cdot x_j, \sum_{j=1}^n \underline{a}_{ij} \cdot x_j, \sum_{j=1}^n \bar{a}_{ij} \cdot x_j \right)_{LR} \end{aligned}$$

where \oplus denotes, as usually, the extended addition operation.

That is to say, in vectorial form:

$$\tilde{y}_i = \tilde{a}_i^t x = (a_i^t x, \underline{a}_i^t x, \bar{a}_i^t x)_{LR}, \text{ being}$$

$$\tilde{a}_i^t = (\underline{a}_{i1}, \underline{a}_{i2}, \dots, \underline{a}_{in}) = (a_i^t, \underline{a}_i^t, \bar{a}_i^t), \text{ where } a_i^t = (a_{i1}, a_{i2}, \dots, a_{in})$$

is the average values vector and

$$\underline{a}_i^t = (\underline{a}_{i1}, \underline{a}_{i2}, \dots, \underline{a}_{in}) \text{ and } \bar{a}_i^t = (\bar{a}_{i1}, \bar{a}_{i2}, \dots, \bar{a}_{in})$$

are the left and right limit vectors respectively.

Note the importance of the formulated hypothesis, since that the addition operation is only well defined if the parts are of the same type (LR, RL, LL or RR) this obliging to consider the variables x_1, \dots, x_n as non-negative.

4.2. Fuzzifying a linear constraints system

From the considerations presented in the previous section, we conclude that the first member of a linear fuzzy constraint is still a fuzzy real number. Consequently, the problem in presence is the definition of (in)equality relation between fuzzy numbers.

Dubois and Prade [8,9] conceive the symbol " $=$ " in two ways:

- (i) as a strict equality between $\tilde{a}_i^t x$ and \tilde{b}_i (equality between membership function), that can be made weaker through the inclusion $\tilde{a}_i^t x \subseteq \tilde{b}_i$.

In this case the diffusion of \tilde{b}_i is interpreted as the maximal tolerance for diffusion of $\tilde{a}_i^t x$, reason why this type of constraints is denominated as tolerance constraints;

- (ii) as an approximate equality between $\tilde{a}_i^t x$ and \tilde{b}_i , by introducing one or more comparison indices between fuzzy numbers.

Therefore, the fuzzy solution \bar{x} of (3.27) is defined by:

$$(3.29) \quad \mathcal{X}(x) = \begin{cases} \alpha & \text{if } x = x^*(\alpha) \\ 0 & \text{otherwise} \end{cases}$$

With regard to $\bar{z}_{op} = \{(r, \mu_z(r)), r \in \mathbb{R}\}$ (def. 3.3), we get

$$\mu_z(r) = \begin{cases} 0 & \text{if } r < 7 \text{ or } r > 16 \\ 1.7 - r/10 & \text{if } 7 \leq r \leq 12 \\ 2 - r/8 & \text{if } 12 \leq r \leq 16 \end{cases}$$

Following this approach, *Verdegay* [32] proposes an example where he considers non-linear membership functions $\mu_{\tilde{C}_j}$.

3.2.2. - Linear Program with punctual solution

Other authors have suggested the determination of a rigid "maximizing decision", associating the o.f., conveniently transformed, with fuzzy constraints of the given problem.

The fundamental question is the circumstance of the domain of the o.f. variation not being normalized, fact that raises difficulties when it is aggregated with the constraints.

- (i) *Tanaka et al* [30] and also *Negoita and Ralescu* [23] started from (3.22) and supposing $f(x) \geq 0, \forall x \in \mathbb{R}^n$, normalised the o.f. using an adequate scale factor.

Effectively, with $M = \sup f(x)$, where $\mathcal{S}(\tilde{C})$ is the support of \tilde{C} and $\mathcal{S}(\tilde{C})$ designates $X \cap \mathcal{S}(\tilde{C})$

it's topological lock.

If $M > 0$, a fuzzy objective may be conceived, if we define $\mu_G(x) = f(x)/M$, verifying $\mu_G(x) \in [0, 1], \forall x \in X \cap \mathcal{S}(\tilde{C})$.

The initial asymmetry is eliminated and in this way the problem can now be solved by the usual methods.

Note that the applied normalisation does not change the eventual linear character of the ordering subjacent to f .

- (ii) *Zimmermann* [41], arguing that the used scale factor does not have, in most cases, a real justification, diffuses the o.f. of (3.22) by associating it with a fuzzy objective whose membership function is linear and defined by:

$$(3.30) \quad \mu_G(x) = \begin{cases} 0 & \text{if } f(x) > \sup_{R_1} f \\ \frac{f(x) - \sup_{R_1} f}{\sup_{\mathcal{S}(R)} f - \sup_{R_1} f} & \text{if } \sup_{R_1} f \leq \sup_{\mathcal{S}(R)} f \\ 1 & \text{if } f(x) \geq \sup_{\mathcal{S}(R)} f \end{cases}$$

where $R = X \cap \tilde{C}$ is the admissible region for problem (3.22), $\mathcal{S}(R) = X \cap \mathcal{S}(\tilde{C})$ and $R_1 = X \cap C_1$ (C_1 is the 1 level cut of the fuzzy set \tilde{C}) and we suppose $\sup_{R_1} f$ and $\sup_{\mathcal{S}(R)} f$ finite quantities.

In this manner, the o.f. being represented by the fuzzy set $\tilde{G} = \{(x, \mu_G(x)), x \in \mathbb{R}^n\}$, with $\mu_G(x) \in [0,1]$, we may recover a totally fuzzy model, where the tecnicas developed in 3.1 are still applicable.

In particular, for the linear model (3.25), and keeping the hypothesis of the previous section, we get:

$$\mathcal{S}(\tilde{C}) = \left(\bigcap_{j=1}^{k_1} \{x \in \mathbb{R}^n : a_j^1 x > b_j - p_j\} \right) \cap \left(\bigcap_{j=k_{1+1}}^{k_2} \{x \in \mathbb{R}^n : a_j^1 x < b_j + p_j\} \right) \cap \left(\bigcap_{j=k_{2+1}}^{k_3} \{x \in \mathbb{R}^n : b_j - p_j < a_j^1 x < b_j + p_j\} \right) \text{ and}$$

$$R_1 = X \cap \left(\bigcap_{j=1}^{k_1} \{x \in R^n : a_j^1 x \geq b_j\} \right) \cap \left(\bigcap_{j=k_1+1}^{k_2} \{x \in R^n : a_j^1 x \leq b_j\} \right) \cap \left(\bigcap_{j=k_2+1}^{k_3} \{x \in R^n : a_j^1 x = b_j\} \right).$$

In this manner, the membership function of the fuzzy objective can be determined by solving the linear programs:

$$\begin{aligned} & \max Z(x) \\ \text{s.a} \quad & a_j^1 x \geq b_j \quad j = 1, \dots, k_1 \\ & a_j^1 x \leq b_j \quad j = k_1 + 1, \dots, k_2 \\ & a_j^1 x = b_j \quad j = k_2 + 1, \dots, k_3 \\ & x \in X \end{aligned}$$

and

$$\begin{aligned} & \max Z(x) \\ \text{s.a} \quad & a_j^1 x \geq b_j - p_j \quad j = 1, \dots, k_1 \\ & a_j^1 x \leq b_j + p_j \quad j = k_1 + 1, \dots, k_2 \\ & a_j^1 x \leq b_j - p_j \quad j = k_2 + 1, \dots, k_3 \\ & a_j^1 x \geq b_j - p_j \quad j = k_2 + 1, \dots, k_3 \\ & x \in X \end{aligned}$$

which lead to $\sup_{R_1} Z(x) = z_1$ and $\sup_{S(R)} Z(x) = z_0$ respectively.

3.3 - Final comments

At this step we have outlined the approaches which we consider more relevant in flexible fuzzy L.P. There are, nevertheless, other important complementary topics that can be used as a background to further studies:

- 1 - We only considered linear membership functions for the representation of the intervenient fuzzy sets, although we have referred other types existing in the literature. Our option lies essentially upon two criteria - the simplicity of manipulation and the frequency with which are applied in the specialised scientific studies.
- 2 - We have limited our work to the models based on the min operator. Nevertheless, we had the opportunity to describe concisely a reasonable range of the available operators. Relatively to it's mathematical applications, we put in relief the *Zimmermann* texts[39] (where the product operator is utilized) and *Luhandjula* [18] (that considers the compensatory operators γ and limited min-sum) in the domain of L.P. with multiple objectives.
- 3 - The models that integrate several objective functions can be studied in the context of some of the described approaches. Namely in relation to the totally fuzzy model, we may find in [11,14,15,16,17,39] an adequate analysis.

As for integer L.P. another interesting extension is fully described in [42].

- 4 - In what concerns the duality of fuzzy L.P. , we point out the work by *Rödder* and *Zimmermann* [26] and subsequently the one of *Hamacher et al* [13], whose point of view is based on the economic interpretation of the dual variables.

Kabbara [15] and *Llena* [17] have conceived a formal study of the duality problem.

Finally, *Verdegay* [32] starts from a new concept of the fuzzy objective to the construction of a dual fuzzy linear program.

On every mentioned approach, it is common the fact that in fuzzy L.P., the pairs of the primal-dual programs admit different formulations (dependent on the adopted model), oppositely the existing uniqueness in classic L.P.

- 5 - As an example of realistic application, we point out the problem formulated and solved by *Wiedey e Zimmermann* [33].

4 - Robust Fuzzy L.P.

In this approach, the difusion of the model is achieved through the constructive coefficients of the o.f. of the problem, and/or through the constraints of the problem, taken as fuzzy parameters mathematically defined by real fuzzy numbers.

Consider the linear program:

Def 3.3 - We designate by "optimal values of a o.f.", T_{op} , the fuzzy set in R , $T_{op} = \{r, \mu_T(r), r \in R\}$, defined by:

$$(3.24) \quad \mu_T(r) = \begin{cases} \sup_{x \in f^{-1}(r)} \mathcal{Z}(x) & \text{if } f^{-1}(r) \neq \emptyset \\ 0 & \text{otherwise,} \end{cases}$$

where $f^{-1}(r) = \{x \in R^n : f(x) = r\}$.

With regard to the linear models, the determination of (\bar{x}, T_{op}) can be made by a parametric linear program [4.31], such as in section 3.1.2. Effectively, let the problem be:

$$(3.25) \quad \max_{x \in C} Z(x) = c^T x,$$

where X as the usual meaning $C = \bigcap_{j=1}^{k_3} \bar{C}_j$, \bar{C}_j being the linear constraint of the type $a_j^T x \geq b_j$ ($j = 1, \dots, k_1$), respectively $a_j^T x \leq b_j$ ($j = k_1+1, \dots, k_2$) and $a_j^T x = b_j$ ($j = k_2+1, \dots, k_3$) characterized by the linear membership function (3.8), (3.9) and (3.10) respectively. We adopt the min operator to define μ_C , i.e., we consider $\mu_C(x) = \bigwedge_{j=1}^{k_3} \mu_{\bar{C}_j}(x)$.

On these conditions, we can write $x_\alpha = x \cap C_\alpha$ on a particular way, the determination of $S(\alpha)$ being made through the solution of a parametric linear program:

$$(3.26) \quad \begin{aligned} & \max Z(x) = c^T x \\ \text{s.a.} & \quad a_j^T x \geq b_j - (1 - \alpha)p_j \quad j = 1, \dots, k_1 \\ & \quad a_j^T x \leq b_j + (1 - \alpha)p_j \quad j = k_1+1, \dots, k_2 \\ & \quad a_j^T x \leq b_j + (1 - \alpha)p_j \quad j = k_2+1, \dots, k_3 \\ & \quad a_j^T x \geq b_j - (1 - \alpha)p_j \quad j = k_2+1, \dots, k_3 \\ & \quad x \in X, \alpha \in [0, 1]. \end{aligned}$$

Let us note that, if T_{op} is a fuzzy set, it is up to the decision maker to choose a pair $(r, \mu_T(r))$ that he consider as optimum, if his aim is to achieved a rigid value.

Ex 3.3 - Let the partially fuzzy linear program, combining only equality type constraints be:

$$(3.27) \quad \begin{aligned} & \max Z(x) = 2x_1 + x_2 \\ \text{s.a.} & \quad x_1 + x_2 \leq 4 \\ & \quad x_1 \leq 3 \\ & \quad x_1 + 2x_2 \leq 6 \\ & \quad x_1, x_2 \geq 0, \end{aligned}$$

relative to which $x = \{(x_1, x_2) \in R^2 : x_1, x_2 \geq 0\}$. Taking $p_1 = p_3 = 4$ and $p_2 = 6$ as the maximal levels of tolerance associated with the intervenient fuzzy constraints, we come up with the following particular version of (3.26).

$$(3.28) \quad \begin{aligned} & \max 2x_1 + x_2 \\ \text{s.a.} & \quad x_1 + x_2 \leq 8 - 4\alpha \\ & \quad x_1 \leq 9 - 6\alpha \\ & \quad x_1 + 2x_2 \leq 10 - 4\alpha \\ & \quad x_1, x_2 \geq 0, \alpha \in [0, 1] \end{aligned}$$

which optimal solution is given by $x^*(\alpha) = \begin{cases} (8-4\alpha, 0) & \text{if } \alpha \leq 1/2 \\ (9-6\alpha, 2\alpha-1) & \text{if } \alpha > 1/2 \end{cases}$

Consequently:

$$S(\alpha) = \{x \in X_\alpha : Z(x) = 16 - 8\alpha\} = \{(8 - 4\alpha, 0)\}, \text{ for } \alpha \leq 1/2$$

and

$$S(\alpha) = \{x \in X_\alpha : Z(x) = 17 - 10\alpha\} = \{(9 - 6\alpha, 6\alpha - 1)\}, \text{ for } \alpha > 1/2$$

4.2.1. Tolerance constraints

Let's consider the constraint $\tilde{a}x \subseteq \tilde{b}$, with $\tilde{a} = (a, \underline{a}, \bar{a})_{LR}$ and $b = (b, \underline{b}, \bar{b})_{LR}$, x being a non negative variable. Then $\tilde{a}x = (ax, \underline{a}x, \bar{a}x)_{LR}$, as previously said. On these conditions the constraint $\tilde{a}x \subseteq \tilde{b}$ is equivalent [8] to the system:

$$\begin{cases} ax = b \\ \underline{a} \leq \underline{b} \\ \bar{a}x \leq \bar{b} \end{cases}$$

More generally, the system of fuzzy linear constraints:

$$\begin{aligned} \tilde{a}_i x &\subseteq \tilde{b}_i & i = 1, \dots, m \\ x &\geq 0, \end{aligned}$$

having the coefficients of LR type, is equivalent to the ordinary system of linear equalities and inequalities:

$$\begin{aligned} \underline{a}_i x &= b_i \\ \underline{a}_i x &\leq \underline{b}_i & i = 1, \dots, m \\ \bar{a}_i x &\leq \bar{b}_i \end{aligned}$$

Ex 4.1 - Let the tolerance constraints system be:

$$\begin{aligned} (4.1) \quad & \tilde{3}x_1 \oplus 4x_2 \subseteq \tilde{13} \\ & \tilde{2}x_1 \oplus (\Theta \tilde{5}x_2) \subseteq \tilde{1} \\ & x_1, x_2 \geq 0, \end{aligned}$$

where $\tilde{3} = (3, 0.2, 0.2)_{RR}$, $\tilde{4} = (4, 0.1, 0.1)_{RR}$, $\tilde{13} = (13, 1, 1)_{RR}$, $\tilde{2} = (2, 0.3, 0.3)_{RR}$, $\tilde{5} = (5, 0.1, 0.1)_{RR}$ and $\tilde{1} = (1, 0.5, 0.5)_{RR}$ are real symmetrical fuzzy numbers whose reference function is defined by $R(x) = \max(0, 1 - x^2)$.

Then we know that $\Theta \tilde{5} = (-5, 0.1, 0.1)_{RR}$ according to the manner used for establishing the symmetric of a fuzzy number of type L.R. (theorem 1.1).

On these conditions, we can write (4.1) in the form:

$$\begin{aligned} (4.2) \quad & \tilde{a}_1 x \subseteq \tilde{b}_1 \\ & \tilde{a}_2 x \subseteq \tilde{b}_2 \\ & x_1 \geq 0, \end{aligned}$$

with $\tilde{a}_1 = (a_1, \underline{a}_1, \bar{a}_1)$ and $\tilde{a}_2 = (a_2, \underline{a}_2, \bar{a}_2)$, $a_1 = (3,4)$, $\underline{a}_1 = \bar{a}_1 = (0.2, 0.1)$ and $a_2 = (2, -5)$, $\underline{a}_2 = \bar{a}_2 = (0.3, 0.1)$.

According to what has been exposed, the given system is equivalent to:

$$(4.3) \quad \begin{aligned} 3x_1 + 4x_2 &= 13 \\ 2x_1 - 5x_2 &= 1 \end{aligned}$$

$$(4.4) \quad \begin{aligned} 0.2x_1 + 0.1x_2 &\leq 1 \\ 0.3x_1 + 0.1x_2 &\leq 0.5 \\ x_1, x_2 &\geq 0, \end{aligned}$$

reduced only to four constraints since the intervening coefficients are symmetric ($\underline{a}_i = \bar{a}_i$ and $\underline{b}_i = \bar{b}_i$ for $i=1,2$).

Analysing the previous conditions, we confirm that the system of the medium values (4.3) admits a single solution $(x_1^*, x_2^*) = (3,1)$ that does not satisfy the second constraint of the system (4.4). In this way, we conclude that the considered fuzzy system does not have a solution.

The approach just described, was inovated by *Dubois* and *Prade* [8,9], and generalizes *Negoita*'s [21] initial proposal, the precussive of the introduction of fuzzy coefficients in linear programs. *Negoita* transformed fuzzy constraints of the type $ax \subseteq \tilde{b}$ in an inlimity of inclusion constraints between the respective α level cuts ($ax_{\alpha} \subseteq b_{\alpha}$, $\alpha \in [0,1]$), and it can be stated that his proposal constitutes, by it's turn, an extension of the model considered by *Soyster* [28].

4.2.2. Approximate equality constraints

This approach is based on a more general problem of comparison between fuzzy real numbers.

In this domain, two types of questions are normally presented: to determinate the lower(higher) fuzzy value of a fuzzy real numbers family and to give a significance to the statement "m̃ is less (greater) than ñ".

For answering the first question, *Dubois* and *Prade* [9] proposed the use of the "m̃n" ("māx") operation defined from the extension principle. In this way given m̃, ñ ∈ $\tilde{\mathcal{N}}(\mathbb{R})$, the fuzzy minimum between m̃ and ñ, designated by m̃n (m̃, ñ) is a fuzzy set in \mathbb{R} , whose membership function is characterized by:

$$\mu_{\min(m, n)}(z) = \sup_{\{(\mu, \nu) \in \mathbb{R}^2 : z = \min(x, y)\}} \min(\mu_{\tilde{m}}(x), \mu_{\tilde{n}}(y))$$

If the supremum in the second member is taken in the set $\{(x, y) : z = \max(x, y)\}$, we obtain the expression for the membership function of māx ($\vee(m, \sim)$, $\vee(n, \sim)$), fuzzy maximum between m̃ and ñ (fig.4.1).

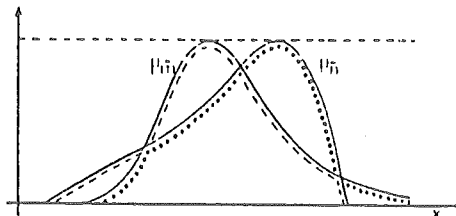


Fig. 4.1. māx (m̃, ñ)
----- m̃n (m̃, ñ)

In [6,9] fundamental results regarding these operations were presented, namely with respect to the properties of the associated approximation formulae for fuzzy numbers of LR type and the corresponding extension for $n(n > 2)$ $\tilde{\mathcal{N}}(\mathbb{R})$ elements.

As for the second question, we find in [3,10] a description of "the major" part of the existing answers, this problem being the focus of studies by several authors. *Dubois* and *Prade* [8,9], generalizing the inequality relation between rigid real numbers, defined a representative index for the "true value" of the statement "m̃ is greater than or equal to ñ", through the extension principle.

$$V(\tilde{m} \leq \tilde{n}) = \sup_{\{(\mu, \nu) \in \mathbb{R}^2 : \mu \geq \nu\}} \min(\mu_{\tilde{m}}(\mu), \mu_{\tilde{n}}(\nu))$$

The following equality is easily verified:

$$\text{alt}(\tilde{m} \cap \tilde{n}) = \min(V(\tilde{n} \leq \tilde{m}), V(\tilde{m} \leq \tilde{n})).$$

It permits to conclude that, when $m \leq n$ (for example), we get $V(\tilde{m} \geq \tilde{n}) = \text{alt}(\tilde{m} \cap \tilde{n})$ since in this case $V(\tilde{n} \geq \tilde{m}) = 1$ (fig. 4.2).

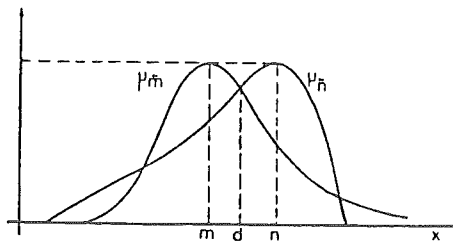


Fig. 4.2 - $V(\tilde{m} \geq \tilde{n}) = \mu_{\tilde{m}}(d) = \mu_{\tilde{n}}(d)$

Nevertheless, the indices $V(\bar{n} \geq \bar{m})$ and $V(\bar{m} \geq \bar{n})$ are simultaneously necessary to make the comparison between \bar{m} and \bar{n} . Effectively, if for example, $V(\bar{m} \geq \bar{n}) = 1$, it is not possible to state that there is a clear separation between \bar{m} and \bar{n} , and it may occur a superposition (even if partial) between \bar{m} and \bar{n} . In this way we are led to:

Def 4.1 - Let $\bar{m}, \bar{n} \in \tilde{N}(R)$ and $\Theta \in [0,1]$. We say that \bar{m} is greater than or equal to n at level Θ , and we write $\bar{m} \geq_{\Theta} \bar{n}$, if $V(\bar{n} \geq \bar{m}) \leq \Theta$ when $m \geq n$.

Def 4.2- Given $\bar{m}, \bar{n} \in \tilde{N}(R)$, we say that m is approximately equal to n at level (ζ, Θ) and we write $\bar{m} \approx_{(\zeta, \Theta)} \bar{n}$, if neither $\bar{m} >_{\zeta} \bar{n}$ nor $\bar{n} \geq_{\Theta} \bar{m}$

Suppose now that $\bar{m} = (m, \alpha, \beta)_{LR}$ and $\bar{n} = (n, \gamma, \delta)_{RL}$.

Then $V(\bar{n} \geq \bar{m}) = L(\frac{m - n}{\alpha + \delta})$ if $m \geq n$ and $V(\bar{m} \geq \bar{n}) = R(\frac{n - m}{\beta + \gamma})$ if $m \leq n$, and we

can write:

$$\begin{cases} \bar{m} >_{\zeta} \bar{n} & \text{iff } m - n \geq \alpha + \delta \quad (\zeta = L(1)) \\ \bar{n} \geq_{\Theta} \bar{m} & \text{iff } n - m \geq \beta + \gamma \quad (\Theta = R(1)), \end{cases}$$

according to the proper of the reference functions L and R.

We can solve a linear fuzzy (in)equality without any difficulty based on the above (in)equality relation between fuzzy real numbers. Suppose that \bar{a} and \bar{b} are fuzzy real numbers of opposite types, $\bar{a} = (a, a_+, \bar{a})_{LR}$ and $\bar{b} = (b, b_+, \bar{b})_{RL}$ respectively and consider the fuzzy equation $\bar{a}x \approx \bar{b}$, with $x \geq 0$. From definition 4.2 and particularizing once again $\Theta = R(1)$, $\zeta = L(1)$, we get the following equivalent form for the given equation:

$$\begin{cases} ax - b < a_+x + \bar{b} & \text{if } ax - b \geq 0 \\ b - ax < b_+ + \bar{a}x & \text{if } ax - b \leq 0 \end{cases}$$

The resolution of the approximated system of equalities:

$$(4.5) \quad \bar{a}_i^+x \approx_{x \geq 0} \bar{b}_i \quad i = 1, \dots, m$$

is now immediate, if we assume that the intervening coefficients of the first member are of LR type and that the independent terms are of RL type. Effectively, in these conditions the system (4.5) takes the form:

$$\bar{y}_i \approx b_i \quad i = 1, \dots, m$$

where $\bar{y}_i = (a_i^+x, a_i^+x, \bar{a}_i^+x)_{LR}$ for $i = 1, \dots, m$. Consequently, by analogy with the unidimensional case, the given system is equivalent to the inequality rigid system:

$$(4.6) \quad \begin{cases} a_i^+x - b_i < a_i^+x + b_i & ; \text{ if } a_i^+x - b_i \geq 0 \\ b_i - a_i^+x < b_i + \bar{a}_i^+x & ; \text{ if } a_i^+x - b_i \leq 0, \quad i=1, \dots, m, x \geq 0 \end{cases}$$

In this way, admitting a rigid o.f., we conclude that from the corresponding fuzzy version, we were led once again to a classic optimization problem, associated with an open polygonal domain.

Ex 4.2 -

Let's consider the fuzzy linear system presented in the previous example, keeping the definition of the coefficients and taking constraints of the approximate equality type (to the adopted levels), that is to say:

$$(4.7) \quad \begin{aligned} \bar{3}x_1 \oplus \bar{4}x_2 &\approx \bar{13} \\ \bar{2}x_1 \oplus (\Theta \bar{3}x_2) &\approx \bar{1} \\ x_1, x_2 &\geq 0 \end{aligned}$$

We can write:

$$\begin{aligned} (a) \quad a_1x - b_1 \geq 0 &\quad \Rightarrow (a_1 - a_1)x < b_1 + b_1 \\ 3x_1 + 4x_2 - 13 \geq 0 &\quad \Rightarrow 2.8x_1 + 3.9x_2 < 14 \\ (b) \quad a_1x - b_1 \leq 0 &\quad \Rightarrow (a_1 + \bar{a}_1)x > b_1 - b_1 \end{aligned}$$

$$3x_1 + 4x_2 - 13 \leq 0 \Rightarrow 3.2x_1 + 4.1x_2 > 12$$

(c) $a_2x - b_2 \geq 0 \Rightarrow (a_2 - \underline{a}_2)x < \underline{b}_2 + b_2$
 $2x_1 - 5x_2 - 1 \geq 0 \Rightarrow 1.7x_1 - 5.1x_2 < 1.5$

(d) $a_2x - b_2 \leq 0 \Rightarrow (a_2 + \bar{a}_2)x > b_2 - \underline{b}_2$
 $2x_1 - 5x_2 - 1 \leq 0 \Rightarrow 2.3x_1 - 4.9x_2 > 0.5$.

The system of strict inequalities is represented in the figure 4.3.

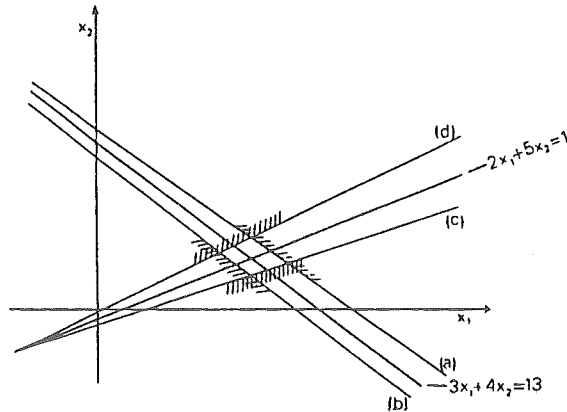


Fig. 4.3 - Solution set of the fuzzy system (4.7)

The developed model is also applicable to systems of fuzzy constraints of the equality type, since the inequality relation presented constitutes a particular case of the introduced equality concept.

It allows also the solution of linear fuzzy programs involving fuzzy objectives, when these take the form of linear constraints through definition of the adequated aspiration levels.

4.3. Extensions

The described approach admits an immediate generalization, if we consider a system of linear constraints where not only the coefficients but also the variables of the problem are fuzzy [8,9]. On this perspective, the presented analysis remains coherent, and makes intervenient the product operation between fuzzy numbers. In particular, with regard to the LR representation the application conditions of the calculation rules stated on theorem 1.2 must be respected when establishing the fuzzy version of the first member of a linear constraint.

On the other hand, we observe that the difusion processes of a linear constraint we have referred are not the unique possibilities. The existing proposals are focused not only on the definition of the constraints 1st member but also on the adoption of other comparison criterion for fuzzy numbers. In this manner, and with respect to the difusion of a linear function, Tanaka and Asai [29] defined the fuzzy set.

$$\tilde{y}_i = \sum_{j=1}^n \tilde{a}_{ij}x_j$$

by application of the extension principle and not using fuzzy real algebra. The same authors proposed the relation "ñ is positive at the level h" ($h \in [0,1)$) that is applied to linear programming.

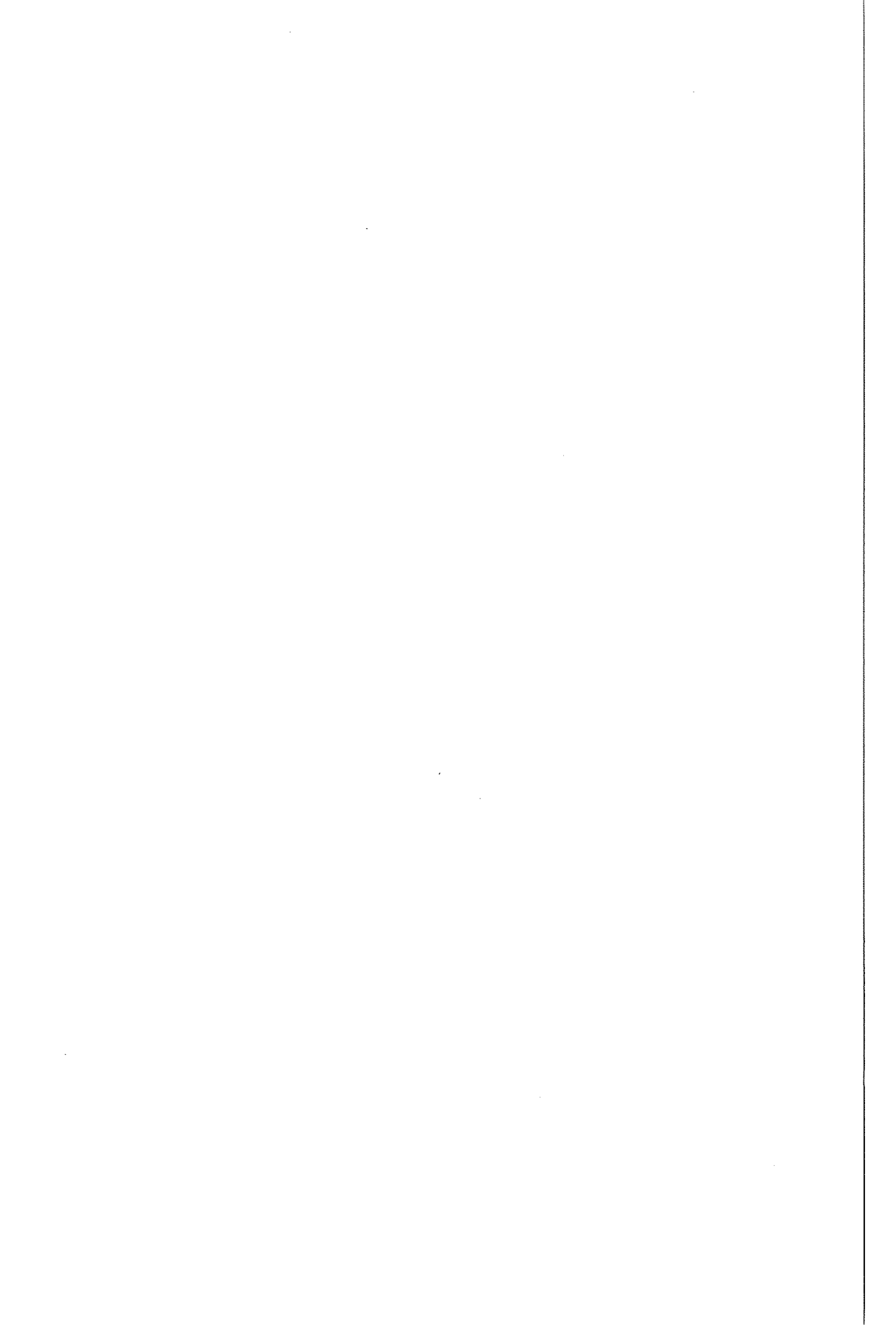
Ramík and Řimánek [25] proposed a inequality concept between elements of $\tilde{N}(\mathbb{R})$ based on the α level cut of a fuzzy set. In the same text, the authors make the considered relation compatible with the maximum and minimum fuzzy notions, indicated in section 4.2.

Dubois [5] established an unified formulation of several approaches for fuzzy L.P. namely the Zimmermann (in flexible programming), Negoita (in robust programming) and Tanaka and Asai [29] (in robust programming) versions which were referred in this work. In that paper, Dubois represents the models coefficients by real fuzzy intervals and uses a set of 4 comparison indices, defined by the author and M. Prade in [10], that are integrated in the theory of possibilities domain, introduced by Zadeh [37] and based on the fuzzy sets theory.

Bibliography

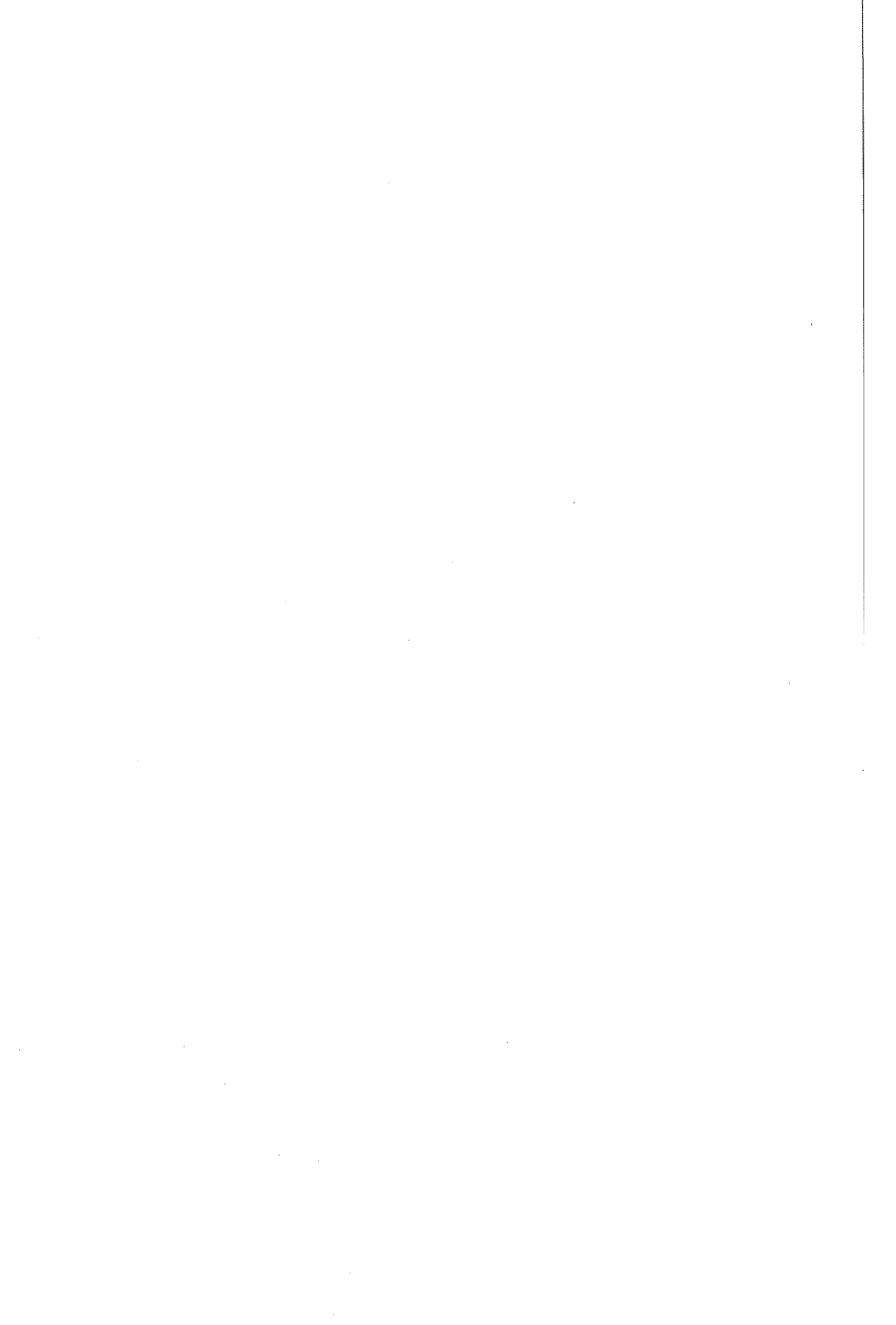
- [1] R. Bellman, M. Giertz - On the analytic formalism of theory of fuzzy sets. *Information Sci.* 5 (1973), 149-156.
- [2] R. Bellman, L. A. Zadeh - Decision making in a fuzzy environment. *Management Sci.* 17(4) (1970), B 141-164.
- [3] G. Bortolan, R. Degani - A review of some methods for ranking fuzzy subsets. *Fuzzy Sets and Systems* 15 (1985), 1-19.
- [4] S. Chanas - The use of parametric programming in fuzzy linear programming. *Fuzzy Sets and Systems* 11 (1983), 243-251.
- [5] D. Dubois - Linear programming with fuzzy data. Being published at : *The Analysis of Fuzzy Information*. J. C. Bezdek, Ed, CRC Press, Boca Raton Fl.
- [6] D. Dubois, H. Prade - Operations on fuzzy numbers. *Internat. J. Systems Sci.* 9(6) (1978), 613-626.
- [7] D. Dubois, H. Prade - Fuzzy real algebra : some results. *Fuzzy Sets and Systems* 2 (1979), 327-348.
- [8] D. Dubois, H. Prade - Systems of linear fuzzy constraints. *Fuzzy Sets and Systems* 3 (1980), 37-48.
- [9] D. Dubois, H. Prade - Fuzzy sets and systems : theory and applications. Academic Press, New York, 1980.
- [10] D. Dubois, H. Prade - Ranking fuzzy numbers in the setting of possibility theory. *Information Sci.* 30 (1983), 83-224.
- [11] R. G. Dyson - Maximin programming, fuzzy linear programming and multicriteria decision making. *J. Opl. Res. Soc.* 31(3) (1980), 263-267.
- [12] J. Flachs, M. A. Pollatschek - Further results on fuzzy mathematical programming. *Information and Control* 38 (1978), 241-257.
- [13] H. Hamacher, H. Leberling, H. J. Zimmermann - Sensitivity analysis in fuzzy linear programming. *Fuzzy Sets and Systems* 1 (1978), 269-281.
- [14] E. L. Hannan - Linear programming with multiple fuzzy goals. *Fuzzy Sets and Systems* 6 (1981), 235-248 .
- [15] G. Kabbara - New utilization of fuzzy optimization method. *Fuzzy Information and Decision Processes* (1982), 239-246, M. Gupta and E. Sanchez, Ed.
- [16] H. Leberling - On finding compromise solutions in multicriteria problems using the min-operator. *Fuzzy Sets and Systems* 6 (1981), 105-118.
- [17] J. Liena - On fuzzy linear programming. *EJOR* 22 (1985), 216-223.
- [18] M. K. Luhandjula - Compensatory operators in fuzzy linear programming with multiple objectives. *Fuzzy Sets and Systems* 8 (1982), 245-252.
- [19] K. G. Murty - Linear programming. John Wiley & Sons, 1983 .
- [20] K. Nakamura - Some extensions of fuzzy linear programming. *Fuzzy Sets and Systems* 14 (1984), 211-229.
- [21] C. V. Negoita - Management applications of system theory. Birkhäuser Verlag, Basel, 1978.
- [22] C. V. Negoita - The current interest in fuzzy optimization. *Fuzzy Sets and Systems* 6 (1981), 216-269.
- [23] C. V. Negoita, D. A. Ralescu - Applications of fuzzy sets to systems analysis. Birkhäuser Verlag, Basel, 1975.
- [24] M. OhEigeartaigh - A fuzzy transportation algorithm. *Fuzzy Sets and Systems* 8 (1982), 235-243.
- [25] J. Ramík, J. Růmánek - Inequality relation between fuzzy numbers and its use in fuzzy optimization. *Fuzzy Sets and Systems* 16 (1985), 123-138.
- [26] W. Rödder, H. J. Zimmermann - Duality in fuzzy linear programming. *Extremal Methods and Systems Analysis* (1980), 415-429, A. Fiacco and K. Kortanek, Ed.
- [27] A. C. Rosa - "Programação Linear Difusa - Uma Análise das Principais Abordagens" - M.Sc. Dissertation, University of Coimbra, 1987.
- [28] A. L. Soyster - Convex programming with set-inclusive constraints and applications to inexact linear programming. *Operations Research* 21 (1973), 1154-1157.
- [29] H. Tanaka, K. Asai - Fuzzy linear programming with fuzzy numbers. *Fuzzy Sets and Systems* 13 (1984), 1-10.
- [30] H. Tanaka, T. Okuda, K. Asai - On fuzzy mathematical programming. *J. Cybernetics* 3(4) (1974), 37-46.
- [31] J. L. Verdegay - Fuzzy mathematical programming. *Fuzzy Information and Decision Processes* (1982), 231-237, M. Gupta and E. Sanchez, Ed.
- [32] J. L. Verdegay - A dual approach to solve the fuzzy linear programming problem.

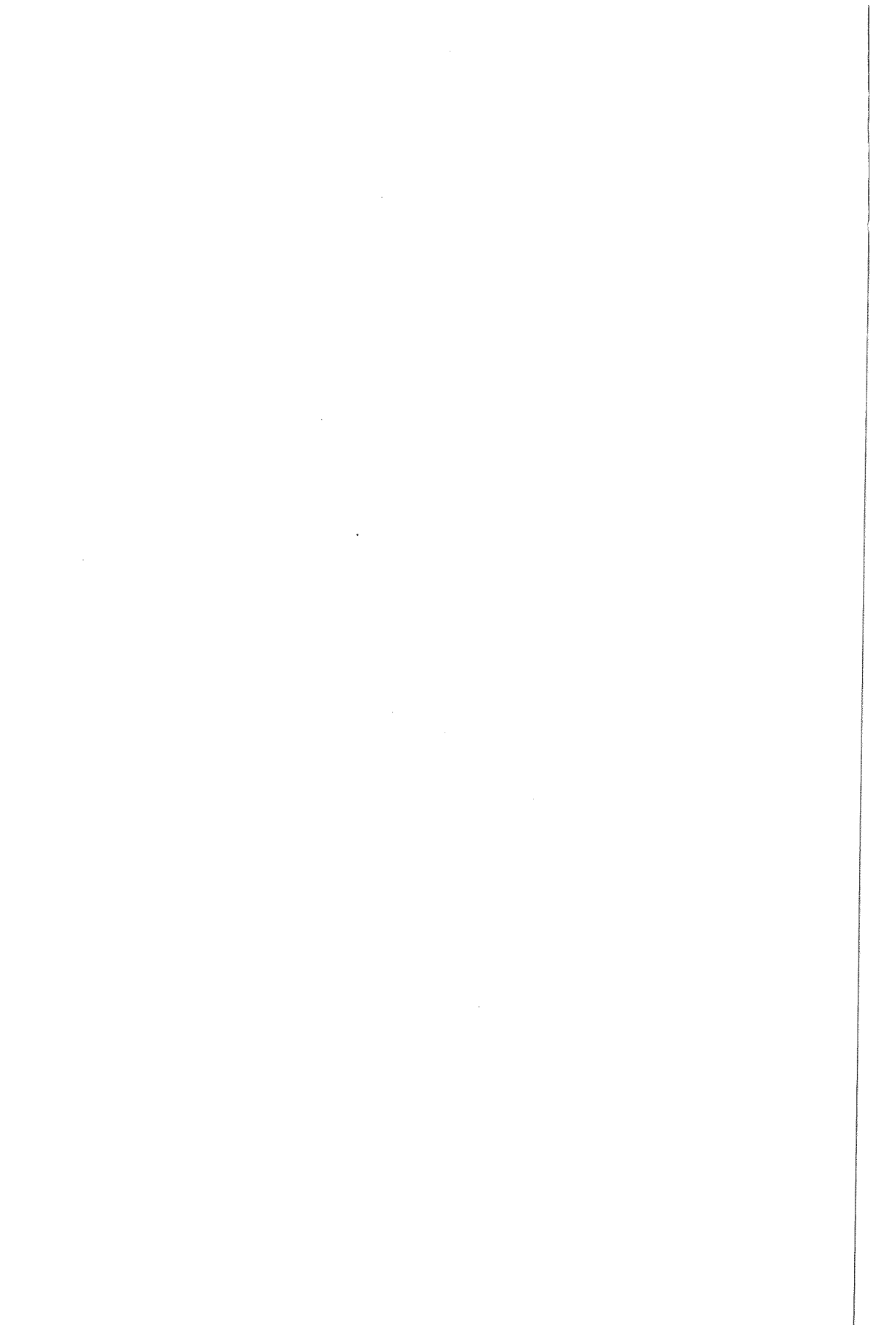
- Fuzzy Sets and Systems 14 (1984), 131-141.
- [33] G. Wiedey, H. J. Zimmermann - Media selection and fuzzy linear programming. *J. Opl. Res. Soc.* 29(11) (1978), 1071-1084.
 - [34] R. Yager - On a general class of fuzzy connectives. *Fuzzy Sets and Systems* 3 (1980), 235-242.
 - [35] L. A. Zadeh - Fuzzy sets. *Information and Control* 8 (1965), 338-353.
 - [36] L. A. Zadeh - The concept of a linguistic variable and its applications to approximate reasoning. Part I : *Information Sci.* 8 (1975), 199-249; Part II : *Information Sci.* 8 (1975), 301-357; Part III : *Information Sci.* 9 (1975), 43-80.
 - [37] L. A. Zadeh - Fuzzy sets as a basis for a theory of possibility. *Fuzzy Sets and Systems* 1 (1978), 3-28.
 - [38] H. J. Zimmermann - Description and optimization of fuzzy systems. *Internat. J. General Systems* 2 (1976), 209-215.
 - [39] H. J. Zimmermann - Fuzzy programming and linear programming with several objective functions. *Fuzzy Sets and Systems* 1 (1978), 45-55.
 - [40] H. J. Zimmermann - Using fuzzy sets in operational research. *EJOR* 13 (1983), 201-216.
 - [41] H. J. Zimmermann - Fuzzy set theory and its applications. *International Series in Management Science/Operations Research*, Kluwer-Nijhoff Publishing, 1985.
 - [42] H. J. Zimmermann, M. A. Pollatschek - Fuzzy 0-1 programs. *Zimmermann et al.* (1984), 133-146.



**FOTOGRAFIA, MONTAGEM,
IMPRESSÃO E ACABAMENTOS**

Secção de Textos da F.C.T.U.C.





ÍNDICE

Maria Teresa Almeida Problema Básico de Distribuição e suas Extensões — Uma Revisão Bibliográfica —	3
<i>Margarida V. Pato, José M. P. Paixão</i> Cutting Planes from Conditional Bounds for Generalized set Covering Problems	13
<i>J. J. Júdice, F. M. Pires</i> Direct Methods for Convex Quadratic Programs Subject to Box Constraints	23
<i>João A. O. Soares, J. A. Assis Lopes</i> Ciclos e Tendência em Séries Económicas: O Pib Português de 1913 a 1986	57
<i>Laira Toscani, Paulo A. S. Veloso</i> Desenvolvimento de Algoritmos Aproximativos por Acercamento: Especificação Formal	65
<i>A. C. Rosa, J. C. N. Clímaco</i> Fuzzy Linear Programming — A Tentative Survey	71



Associação Portuguesa para o Desenvolvimento
da Investigação Operacional

**CÉSUR — Instituto Superior Técnico — Avenida Rovisco Pais
1000 Lisboa — Telef. 86 74 55**